

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 : 29. 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, 29. April 2019

Autor Name

Inhaltsverzeichnis

Kurzfassung	V
Danksagung	VI
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	5
2.1 Unterscheidung von Daten und Informationen	5
2.1.1 Personenbezogene Daten	5
2.2 Social Engineering	6
2.2.1 Phishing	7
2.2.2 Spear-Phishing	7
2.3 Open Source Intelligence	8
2.4 Web Crawler	8
2.5 Web Scraper	9
2.5.1 Natural Language Processing	9
3 Problembeschreibung	11
4 Ethische und rechtliche Betrachtung	12
5 Anforderungsanalyse	13
5.1 Anforderung an das OSINT einer ausgewählten Person	13
5.2 Anforderung an die Generierung einer Phishing-Mail	14
5.2.1 Anforderung an die Generierung der E-Mail-Adressen	14
5.2.2 Anforderung an die Erstellung der E-Mail-Texte	14

6	Lösungsideen	15
6.1	Methoden für das OSINT einer ausgewählten Person	15
6.1.1	Verwendung von OSINT-Tools	15
6.1.2	Entwicklung eines OSINT-Algorithmus	15
6.2	Konzept für die Erstellung einer Phishing-Mail	16
6.2.1	Methoden zur Generierung der E-Mail-Adresse	16
6.2.2	Methoden zur Generierung des E-Mail-Textes	17
7	Bewertung der Lösungsideen anhand den Anforderungen	19
7.1	Bewertung der OSINT-Methoden für eine ausgewählte Person	19
7.2	Bewertung der Methoden zur Erstellung einer Phishing-Mail	20
7.2.1	Generierung der E-Mail-Adressen	20
7.2.2	Generierung der E-Mail-Texte	20
8	Implementierung des OSINTs für eine ausgewählte Person	22
8.1	Auswahl der Programmiersprache	22
8.1.1	Ziele und Anforderungen	22
8.1.2	Lösungsideen	22
8.1.3	Bewertung der Lösungsideen anhand den Anforderungen	23
8.2	Personensuche im Internet	23
8.2.1	Ziele und Anforderungen	23
8.2.2	Lösungsideen	23
8.2.3	Bewertung der Lösungsideen anhand den Anforderungen	25
8.3	Umsetzung der Personensuche mit Hilfe der Google-Suchmaschine	26
8.3.1	Eingabe der bekannten Daten	26
8.3.2	Generierung der Google-Such-URLs	28
8.3.3	Auswahl der Bibliothek für Serveranfragen	31
8.3.4	Erstellung des Web Crawlers	32
8.4	Die Personenidentifizierung	38
8.4.1	Ziele und Anforderungen	38
8.4.2	Lösungsideen	38
8.4.3	Bewertung der Lösungsideen anhand den Anforderungen	40
8.5	Umsetzung der Methode zur Verwendung eines Identifikationsschlüssels	40
8.6	Herausfiltern von wichtigen Informationen auf einer Webseite	41
8.6.1	Ziele und Anforderung	41
8.6.2	Lösungsideen	41
8.6.3	Bewertung der Lösungsideen anhand den Anforderungen	42
8.7	Umsetzung der Methoden zum Herausfiltern von wichtigen Informationen	43
8.7.1	Text formatieren	43
8.7.2	Erstellung der Wortsammlungen	43
8.7.3	Automatisierte Schlüsselwortgenerierung	44

8.7.4	Umsetzung des Automatic Keyword Extraction mit NLP	48
8.7.5	Herausfiltern von Geburtsjahren	49
8.7.6	Herausfiltern von E-Mail-Adressen	50
8.8	Auswahl der gefundenen Information	51
8.8.1	Ziele und Anforderungen	51
8.8.2	Lösungsidee	52
8.8.3	Bewertung der Lösungsidee anhand den Anforderungen	53
8.9	Umsetzung der Methode zur Auswahl der gefundenen Information	54
8.10	Kontaktanalyse	55
8.10.1	Ziele und Anforderungen	55
8.10.2	Lösungsideen	56
8.10.3	Bewertung der Lösungsideen anhand den Anforderungen	56
8.11	Umsetzung der Instagram-Kontaktanalyse	56
8.11.1	Auslesen der Kontakte	57
8.12	Speicherung der gewonnenen Daten	59
8.12.1	Ziele und Anforderungen	59
8.12.2	Lösungsideen	60
8.12.3	Bewertung der Lösungsideen anhand den Anforderungen	60
8.13	Umsetzung der Personenklasse	60
9	Implementierung der Phishing-E-Mail	62
9.1	Umsetzung des Algorithmus zur Generierung von E-Mail-Adressen	62
9.1.1	Funktion des Algorithmus	62
9.2	Umsetzung der E-Mail-Muster	65
9.2.1	Kategorisierung der E-Mail-Muster	65
9.3	Versenden einer Phishing-E-Mail	70
9.3.1	Ziele und Anforderungen	70
9.3.2	Lösungsideen	70
9.3.3	Bewertung der Lösungsideen anhand den Anforderungen	70
9.4	Umsetzung zum Versenden einer Phishing-E-Mail	71
10	Evaluation der Implementation	72
10.1	Validierung des Gesamtkonzeptes	72
10.2	Beschreibung und Motivation der Testfälle	73
10.2.1	Testfall 1	73
10.2.2	Testfall 2	74
10.2.3	Testfall 3	75
10.3	Übersicht und Bewertung der erzielten Ergebnisse	77
10.3.1	Bewertung von Testfall 1	77
10.3.2	Bewertung von Testfall 2	78
10.3.3	Bewertung von Testfall 3	78

11 Fazit und Ausblick	79
11.1 Fazit	79
11.2 Ausblick	80
A Ausschnitt der Wortsammlung “Tätigkeit”	82
Literatur	84

Kurzfassung

In dieser Abschlussarbeit wird eine Open-Source-Intelligence-Anwendung (OSINT-Anwendung) entwickelt, welche automatisiert personenbezogene Daten zu einer angegebenen Person herausucht und die gefundenen Daten selbständig in eine Phishing-E-Mail integriert. Es wird mit einem kritischen Blick aufgezeigt, wie veröffentlichte personenbezogenen Daten für einen Social-Engineering-Angriff missbraucht werden können. Dabei stellen sich die folgenden Fragen: Ist es möglich, ein Opferprofil aus den gewonnenen Daten zu erstellen, welches ausschließlich korrekte Informationen enthält? Mit welchem Aufwand ist ein automatisierter Spear-Phishing-Mail-Angriff verbunden?

Hierfür wurden Methoden entwickelt um eine Person weitestgehend zu identifizieren, wichtige Informationen aus Webseiten herauszulesen, Phishing-Mail-Muster zu erzeugen und E-Mail-Adressen aus den gewonnenen Daten zu generieren. Die Ergebnisse dazu zeigen, dass mit den erstellten Methoden kein vollständig korrektes Opferprofil erstellt werden kann. Dennoch ist die Generierung einer korrekten Phishing-Mail unter Verwendung der gefundenen Daten möglich.

Somit besteht die Gefahr einer missbräuchlichen Verwendung von personenbezogenen Daten durch eine automatisierte OSINT-Anwendung, für die ein geringer Aufwand benötigt wird.

Danksagung

An dieser Stelle möchte ich mich bei all denjenigen bedanken, die mich bei meiner Bachelorarbeit und während meines Studiums unterstützt haben.

Mein Dank gilt meinen Eltern Wolfgang und Birgit Lang ohne die mein Studium in dieser Weise nicht möglich gewesen wäre.

Ebenfalls möchte ich meinem Vater, meinem Bruder und meiner Freundin für das Korrekturlesen dieser Bachelorarbeit danken.

Ein besonderer Dank geht an meinen Prüfer Prof. Dr. rer. nat. Eggendorfer für die zahlreichen Anregungen und Tipps, sowie für das Korrektur lesen.

1 Einleitung

1.1 Motivation

Das Internet ist die meistgenutzte Informationsquelle und birgt dadurch Gefahren für jeden einzelnen Anwender. [All18] Trotzdem sehen sich laut einer Umfrage nur 30 % der Internetnutzer durch das Risiko einer missbräuchlichen Verwendung ihrer Daten nach einem Hack gefährdet. [Ang18] Jedoch zeigen aktuelle Beispiele, dass jede Person von Datendiebstahl im Netz bedroht ist: Die Schadsoftware Emotet nutzt bereits jetzt die Informationen aus dem sogenannten “Outlook-Harvesting“. Hierbei werden gefälschte E-Mails mit Schadsoftware im Namen von Freunden, Kollegen oder Nachbarn versendet. [Bun]

Jedoch muss nicht zwingend ein Hacker-Angriff einer solchen Gefährdung vorausgehen. Auch freiwillig geteilte Informationen, beispielsweise auf Social-Media-Plattformen wie Facebook, Instagram und LinkedIn, können missbraucht werden. Diese Informationen sind personenbezogen und können deshalb bei einem Phishing-Mail-Angriff genutzt werden, um ein potentielles Opfer zu täuschen oder zu manipulieren.

Solch ein Angriff benötigt allgemein im Voraus Informationen über das Opfer. Neben den Daten aus Social-Media-Plattformen können dabei auch beliebige Einträge im Internet als Quelle dienen, um das Opferprofil zu erweitern. Durch kostenlose und frei zugängliche Open-Source-Intelligence-Tools (OSINT-Tools), welche Informationen über Opfer von öffentlichen Medien sammeln, ist dies für jeden Internetnutzer ohne große Vorkenntnisse möglich.

Dennoch ist ein Phishing-Mail-Angriff unter Verwendung dieser OSINT-Tools mit einem großen Aufwand verbunden. Das hat den Grund, dass die einzelnen Schritte wie die Suche

nach Informationen, Datenanalyse und Integrierung der Daten in eine Phishing-Mail nicht vollautomatisch ablaufen. Die Untersuchung zur Automatisierung dieser Schritte ist notwendig, um die Gefahr der missbräuchlichen Verwendung von öffentlich frei zugänglichen personenbezogenen Daten aufzuzeigen.

1.2 Zielsetzung und Forschungsfragen

Ziel ist es eine OSINT-Anwendung zu entwickeln, welche automatisiert nach personenbezogenen Daten im Internet sucht. Die gefundenen 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:

Ist es möglich ein Personenprofil zu erstellen, bei dem ausschließlich korrekte Informationen vorhanden sind?

Mit welchem Aufwand ist ein automatisierter Spear-Phishing-Mail-Angriff verbunden?

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

Die zu erstellende 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. Die Herausforderung besteht darin, zu erkennen, wann es sich um eine Information der gesuchten Person handelt.

Ziel 2: E-Mail-Adressen finden oder aus den gewonnen Daten generieren.

Wenn eine E-Mail-Adresse zu einer gesuchten Person nicht gefunden werden kann, soll diese mit Hilfe der gewonnen 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. Das Ziel ist in diesem Fall, 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 heraussucht. 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 den 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 des Umgangs mit personenbezogenen Daten aufgezeigt, auf welches in dieser Arbeit eingegangen wird. Darauf folgt die ethische und rechtliche Betrachtung in Kapitel 4. Als nächstes werden die Anforderungen im Kapitel 5 festgelegt und analysiert. Darauf folgen die Lösungsvorschläge im Kapitel 6 und die Auswahl der Lösung anhand den Anforderungen aus Kapitel 7. Anschließend wird bei der Umsetzung auf den praktischen

Teil eingegangen. Dieser unterteilt sich in die Themen Implementierung des OSINTs für eine ausgewählte Person in Kapitel 8 und die Implementierung der Phishing-Mail in Kapitel 9. Am Ende dieser Arbeit befindet sich die Evaluation der Implementation in Kapitel 10 sowie das Fazit und der Ausblick in Kapitel 11.

2 Grundlagen

2.1 Unterscheidung von Daten und Informationen

Als Daten werden (Zahlen)-Werte, Angaben, und formulierbare Befunde bezeichnet, welche unter anderem durch Beobachtungen, Messungen oder statistische Erhebungen gewonnen werden. [Gmb18] Umgangssprachlich wird von Gegebenheiten, Tatsachen und Ereignissen gesprochen. [Pro15]

Um aus Daten Informationen zu erzeugen, wird diesen eine Bedeutung zugeordnet. [Pro15] Das heißt, eine beispielhafte Ziffernfolge “12345“ kann mit Hilfe eines Kontexts zu einer Postleitzahl oder einer Telefonnummer führen.

2.1.1 Personenbezogene Daten

Die DSGVO definiert **personenbezogene Daten** in, 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. [Art.4 Nr.1 DSGVO]

2.2 Social Engineering

Der Grundgedanke von Social Engineering ist, eine Zielperson so zu manipulieren, dass sie für den Angreifer eine bessere Entscheidung trifft. [Had11]

Kevin D. Mitnick definiert Social Engineering in die Überzeugungskraft, Menschen zu täuschen, indem der Social Engineer ein Opfer manipuliert oder vorgibt jemand zu sein, welcher er nicht ist. Dadurch kann er Menschen ausnutzen, um an Informationen zu gelangen. Dabei spielt es keine Rolle, ob Technologien verwendet werden oder nicht. [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. Eine bekannte Angriffsmethode dafür ist das Phishing, auf welche in dieser Arbeit detailliert eingegangen wird.

Der Aufbau eines Social-Engineering-Angriffes ist in mehrere Phasen definiert. 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 eigentliche 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. [DAT15]

2.2.1 Phishing

Das Wort Phishing wird von dem Wort “fishing“ abgeleitet. Der Vergleich ist zutreffend, da der Angreifer ähnlich wie ein Angler am See nach Informationen “fischt“. 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 ist, dass der Angreifer eine gefälschte E-Mail von von dem Onlineversandhaus “Amazon.com, Inc.“ 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.

Dieses Verfahren ermöglicht Angreifern die Anmeldedaten von einer Zielperson sich 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 eine gezielte Mail an ein ausgewähltes Opfer 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.

2.3 Open Source Intelligence

Open Source Intelligence (OSINT) ist definiert als alle Informationen, welche aus öffentlich frei zugänglichen Quellen gewonnen werden. Dabei kann sich die Bedeutung fallspezifisch ändern. So bedeutet OSINT für die CIA die Informationsgewinnung aus ausländischen Nachrichtensendungen. Grundsätzlich ist mit dem Begriff OSINT jedoch die Gewinnung eines öffentlichen Inhalts aus dem Internet gemeint. Eine Verbindung mit dem Begriff Open-Source-Software besteht nicht. [Baz18]

2.4 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.

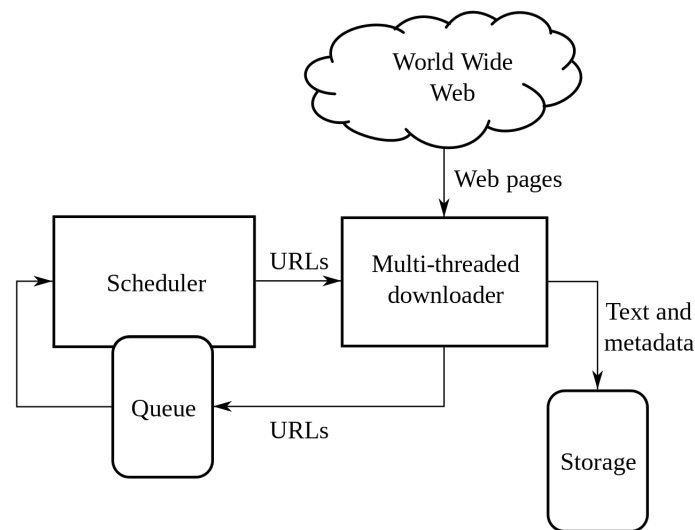


Bild 2.1: Architektur eines Web Crawlers [Cas08]

2.5 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]

2.5.1 Natural Language Processing

Natural Language Processing (NLP) beschreibt eine Technologie für die Kommunikation zwischen Mensch und Computer. Dabei ist das 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 und die Generierung von N-Grammen. [Lit16]

2.5.1.1 Stemming

Stemming ist eine Methode der Wortstandardisierung, bei der verwandte Wörter auf ihrer Stammform reduziert werden. Dabei wird bei dem Rechenvorgang 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 “Kampf“, “kämpfen“ und “kampflös“, welche alle auf den Stamm “kampf“ reduziert werden könnten. Die Verwendung von Stemming kann bei der Schlüsselwortgenerierung von Texten sehr hilfreich sein, da die Anzahl der möglichen Schlüsselwörter reduziert werden können.

2.5.1.2 N-Gramme

N-Gramme sind das Resultat einer Zeichenkettenzerlegung. Hierfür wird ein String in Fragmente unterteilt. Dabei zählen N zusammenhängende Fragmente als ein N-Gramm. [CT⁺94]

In dieser Arbeit besteht ein Fragment immer aus einem vollständigen Wort. Das folgende Beispiel soll dies noch einmal verdeutlichen:

Die Zeichenkette “Das ist ein Beispiel“ wird in die Fragmente [’Das’], [’ist’], [’ein’] und [’Beispiel’] aufgeteilt. Dadurch entstehen die Monogramme mit dem Wert 1 für N. Weiterführend können aus der Zeichenkette die Bigramme [’Das ist’], [’ist ein’], [’ein Beispiel’] erstellt werden, indem N die Zahl 2 beträgt. Identisch zu den vorigen Schritten können Trigramme, Tetragramme, et cetera erstellt werden. Dabei darf die Zahl N nicht die Anzahl der Wörter in der Zeichenkette übersteigen.

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 Social-Media-Seiten wie Twitter, Facebook und Instagram. Allerdings wird auch auf anderen Webseiten personenbezogene Daten in großen Mengen bereitgestellt. Ein Beispiel dafür ist das 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.

Das hier beschriebene Problem zeigt, dass der Zugang für persönliche Information durch das Internet für die Öffentlichkeit einfacher gemacht wird. Es wird mit einem kritischen Blick aufgezeigt, 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 gelangen sollten. Des Weiteren ist der hier verwendete Crawler nicht stark genug, um die Leistung eines Servers von einem sozialen Netzwerk zu beeinflussen. Darüber hinaus werden die gefundenen personenbezogenen Daten nicht gespeichert.

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 Realität durchgeführt wird.

5 Anforderungsanalyse

Die im Kapitel 1.2 definierten Ziele sollen mit den folgenden Anforderungen gewährleistet werden. Dabei unterteilt sich dieser Abschnitt in die Anforderungen an das OSINT einer ausgewählten Person und an die Generierung einer Phishing-Mail.

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, Geschlecht, Wohnort beziehungsweise Standort, 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 werden.

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 großen Pool an generierten E-Mail-Adressen soll die Wahrscheinlichkeit erhöht werden, dass die richtige E-Mail-Adresse dabei ist. Darüber hinaus können die Adressen 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. Diese sollen die gewonnenen Informationen zur Generierung verwenden, damit für jedes Opfer ein übereinstimmender Text erstellt werden kann. Die Texte sollen mit den gefundenen Daten Sinn ergeben und eine korrekte Grammatik beinhalten. Weiterführend müssen Social-Engineering-Fähigkeiten genutzt werden, um menschliche Wünsche, Ängste und verbreitete Verhaltensmuster der Zielperson auszunutzen. [DAT15] Dadurch soll von dem möglichen Misstrauen des Opfer gegenüber der Phishing-Mail abgelenkt werden.

6 Lösungsideen

In diesem Kapitel werden die Lösungsideen für die Umsetzung der im Kapitel 1.2 definierten Ziele beschrieben. Zu Beginn werden Methoden für das OSINT einer ausgewählten Person vorgestellt. Anschließend wird ein Konzept zur Erstellung einer Phishing-E-Mail aufgezeigt.

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 bereitstellt, ist die von Michael Bazzell [Baz19b]. Sie stellt Methoden zur Suche nach E-Mail-Adressen, Benutzernamen, Social-Media-Profilen, et cetera zu Verfügung. Allerdings werden nicht nur selbst entwickelte OSINT-Methoden von Michael Bazzell bereitgestellt, sondern auch andere Webseiten mit weiteren OSINT-Tools vorgeschlagen.

6.1.2 Entwicklung eines OSINT-Algorithmus

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 einen 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.

6.2.1.1 Entwicklung eines Algorithmus zur Adressgenerierung

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 ein Pool mit möglichen Mailadressen erstellt. Dabei muss jede einzelne E-Mail-Adresse auf Validität geprüft werden.

6.2.1.2 Verwendung von automatisierbaren OSINT-Tools

Für die Generierung der E-Mail-Adressen kann ein kostenloses OSINT-Tool von Michael Bazzel verwendet werden. Dieses Tool ermöglicht es, die gewonnenen Informationen über ein Formular einzugeben. Anschließend werden daraus mögliche E-Mail-Adressen generiert. Auch hier entsteht ein Adresspool, bei dem die E-Mail-Adressen auf Validität geprüft werden. Zu dem hat das Tool eine weitere Funktion. Es wird automatisiert nach Einträgen der generierten E-Mail-Adressen im Internet gesucht und angezeigt. [Baz19a]

6.2.2 Methoden zur Generierung des E-Mail-Textes

6.2.2.1 Erstellung von E-Mail-Mustern

Die zu erstellenden E-Mail-Muster entsprechen hier kategorisierten Lückentexten, welche frei erstellt werden. 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 wäre in die Kategorien Privat und Geschäftlich. Dabei werden bei allen Texten menschliche Gefühle wie Freude und Ängste angesprochen. Beispielsweise können bei einem Student Ängste hervorgerufen werden, indem die E-Mail auf einen verpassten Rückmeldezeitraum hinweist.

6.2.2.2 Erzeugung von Text-Fragmenten

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 6.2.2.1 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 bestehen. Diese Anzahl kann variieren, da sie abhängig von den gefundenen Informationen über die Zielperson ist.

7 Bewertung der Lösungsideen anhand den Anforderungen

7.1 Bewertung der OSINT-Methoden für eine ausgewählte Person

Es gibt zwei verschiedene Methoden um OSINT zu betreiben. Die erste Lösungsidee beschreibt die Verwendung von einem öffentlich frei zugänglichen OSINT-Tool. Dieses Tool bietet zahlreiche Möglichkeiten um eine Person beziehungsweise Daten über eine Person zu finden. Allerdings ist es auf dieser Webseite nicht möglich, ein Profil zur gesuchten Person anzugeben. Es kann nur eine begrenzte Anzahl an Daten über ein Formular eingegeben werden. Des Weiteren wird bei einer Suche ausschließlich nach dem Namen oder einer E-Mail gesucht. Das Suchergebnis ist dadurch kein vollständiges Personenprofil, es werden lediglich Verweise auf weitere Webseiten mit möglichen Einträgen angezeigt.

Im Gegensatz zu diesem Tool nutzt der eigene Algorithmus alle im Vorfeld bekannten Daten für eine Suche. Es wird nicht auf Webseiten mit möglichen Informationen verwiesen, sondern Profile der Zielpersonen erstellt. Durch die Verwendung eines eigenen Algorithmus kann die Suche an die Anforderungen beliebig angepasst werden. Zusätzlich besteht die Möglichkeit zur Optimierung der Suche durch die Verwendung der Techniken aus [Baz18].

7.2 Bewertung der Methoden zur Erstellung einer Phishing-Mail

7.2.1 Generierung der E-Mail-Adressen

Ein bereitgestelltes OSINT-Tool ist ein komplettes und funktionsfähiges System. Dadurch wird kein zusätzlicher Aufwand für die Entwicklung eines Algorithmus benötigt. Lediglich die Automatisierung des Tools muss erstellt werden. Allerdings kann nicht jede Information über die Person zur Generierung der E-Mail-Adresse genutzt werden. Dies ist ein großer Nachteil. Die Wahrscheinlichkeit, dass sich die richtige Adresse unter den erzeugten befindet, sinkt dadurch.

Bei einem eigenen Algorithmus fließen dagegen alle Information mit in die Generierung einer E-Mail-Adresse ein. Beispielsweise auch das Geburtsjahr einer Zielperson, was bei dem OSINT-Tool [Baz19a] nicht verwendet wird. Jedoch können die vom OSINT-Tool generierten Adressen, als Anregung und Ideengeber für den eigenen Algorithmus dienen. Für die erfolgreiche Simulation eines Phishing-Mail-Angriffes wird die korrekte E-Mail-Adresse benötigt. Aus diesem Grund wird der eigene Algorithmus verwendet. Dadurch wird die Wahrscheinlichkeit erhöht, dass sich die korrekte Mail-Adresse in dem Pool befindet.

7.2.2 Generierung der E-Mail-Texte

Der Inhalt einer E-Mail ist sehr wichtig für die Glaubwürdigkeit einer Phishing-Mail. Es ist die einzige Möglichkeit, bei der ein Angreifer mit dem Opfer kommunizieren kann. Aus diesem Grund ist es von Bedeutung, dass der E-Mail-Text Sinn ergibt und eine korrekte Grammatik enthält. Des Weiteren können dabei Gefühle und Verhaltensmuster des Opfers ausgenutzt werden.

Bei der dynamischen Texterzeugung mit der Fragment-basierten Methode 6.2.2.2 kann die Grammatik und der Zusammenhang des Textes zu einer Problematik führen. Durch die Verkettung von verschiedensten Fragmenten könnte ein Text erzeugt werden, welcher keinen sinnvollen 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. Das hat den Grund, dass ein Opfer dadurch vermutet, dass dieser Text von einem Menschen verfasst wurde. Somit werden die vollständigen E-Mail-Muster umgesetzt.

8 Implementierung des OSINTs für eine ausgewählte Person

8.1 Auswahl der Programmiersprache

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

8.1.1 Ziele und Anforderungen

Es 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.

8.1.2 Lösungsideen

Für eine web-basierende Anwendung eignet sich eine dynamische Programmsprache, da ein Programm zur Laufzeit erweitert werden kann. Infolgedessen zählen Python und Ruby als mögliche Programmiersprache.

8.1.3 Bewertung der Lösungsideen anhand den Anforderungen

Beide Sprachen können Webseiten, welche JavaScript-Code enthalten, 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 die vorgegebenen Ziele umzusetzen. Somit haben beide Programmiersprachen die Voraussetzungen für die Entwicklung der Anwendung. Allerdings bietet Python in diesem Bereich eine größere 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 Personensuche im Internet

8.2.1 Ziele und Anforderungen

Ziel ist es Methoden zu entwickeln, welche eine Zielperson auf dem optimalen Weg im Internet suchen. Dazu soll die Suche mit Hilfe der eingegebenen Daten eingegrenzt werden, dass möglichst viele zutreffende Informationen über die gesuchte Person gefunden werden kann. Gleichzeitig sollen Fehlergebnisse vermieden werden. Allerdings dürfen dabei nicht zu viele Informationsquellen ignoriert werden, da sonst wichtige Information verloren gehen könnten.

8.2.2 Lösungsideen

Die Art der Personensuche wird abhängig von den eingegebenen Daten variiert. Das heißt, dass die eingegebenen Daten über die Zielperson vor der Suche analysiert werden. Abhängig von den Ergebnissen der Analyse wird die Personensuche durchgeführt. Für die Implementierung der Suche werden nachfolgend zwei mögliche Methoden beschrieben.

8.2.2.1 Personensuche mit Hilfe einer Suchmaschine

Bei dieser Methode wird mit Hilfe einer Suchmaschine nach Informationen gesucht. Mögliche Suchmaschinen sind “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 weiteres Szenario beschreibt der Fall, 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.

8.2.2.2 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.

8.2.3 Bewertung der Lösungsideen anhand den Anforderungen

Um möglichst viele zutreffende Informationen über eine Person im Internet zu finden, bietet die Personensuche mit der Verwendung einer Suchmaschine die beste Lösung. Das hat den Grund, dass das ganze Internet nach Informationen durchsucht wird, anstatt ausschließlich auf festgelegten Webseiten. 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 Suche mit einer Suchmaschine zusätzliche Informationen herausgefunden werden, die für das spätere OSINT von Vorteil sind. Falls sich eine Social-Media-Seite unter den Suchergebnissen befindet, kann diese nachträglich ebenfalls durchsucht werden.

8.2.3.1 Auswahl der Suchmaschine

8.2.3.1.1 Ziele und Anforderungen

Es soll eine Suchmaschine ausgewählt werden, welche einen großen Marktanteil in Deutschland erweist. Dadurch sollen konkrete Suchanfragen möglich sein.

8.2.3.1.2 Lösungsideen

Eine Möglichkeit zur Auswahl ist die Suchmaschine Google von "Google LLC". Eine Alternative hierfür stellt die Suchmaschine Bing dar.

8.2.3.1.3 Bewertung der Lösungsideen anhand den Anforderungen

Laut einer Expertenaussage sucht Bing tiefgreifender nach Information auf Social-Media-Plattformen wie Facebook, Twitter und LinkedIn. [Boh14] Allerdings finden nur 3,5% aller Suchanfragen in Deutschland über Bing statt. Im Gegensatz dazu hat Google einen Marktanteil von 91,2% in Deutschland. [Boh14]

Diese Zahlen sprechen eindeutig für Google. Durch die höhere Anzahl von Suchanfragen, können mehr Daten erfasst und die Ergebnislisten besser bewertet werden. Dies hat zu Folge, dass Bing bei einer konkreten Suche schlechter abschneidet.

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 Vorteil für Bing. Das hat den Grund, dass bei der verwendeten Suche standardmäßig die bekannten Social-Media-Plattformen nach Daten kontrolliert werden.

8.3 Umsetzung der Personensuche mit Hilfe der Google-Suchmaschine

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

8.3.1 Eingabe der bekannten Daten

Es besteht die Möglichkeit den **Vorname**, **Nachname**, **Geburtsjahr**, **Geschlecht**, **Wohnort** bzw. **Standort** , **Arbeitgeber**, **Instagram-Benutzername**, **Facebook-Benutzername**, **Twitter-Benutzername** und die **E-Mail-Adresse** der gesuchten Person über eine Konsole einzugeben.

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

8.3.1.1 Verarbeitung der Daten

Im ersten Schritt wird kontrolliert, welche Informationen von dem Programm-Nutzer eingegeben wurden. Der Vorname und Nachname sind nicht ausreichend für die Suche. Es wird mindestens ein weiteres Attribut benötigt. Dagegen ist der Benutzername 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. Das hat den Grund, dass es für die Personenerkennung wichtig ist, 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/Standort;

Vorname, Nachname, Geburtsjahr;

Vorname, Nachname, Institution;

Vorname, Nachname, Wohnort/Standort, Geburtsjahr;

Vorname, Nachname, Wohnort/Standort, 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 zu jeder Information ein Eintrag im Internet besteht.

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

8.3.2 Generierung der Google-Such-URLs

8.3.2.1 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]

$\langle scheme \rangle : \langle scheme - specific - part \rangle$ [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]

$http : // \langle host \rangle : \langle port \rangle / \langle path \rangle ? \langle searchpart \rangle$ [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 verändern sich.

Für den $\langle host \rangle$ kann der Fully Qualified Domain Name (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 $\langle path \rangle$ stellt ein HTTP-Selektor dar und ist mit einem Fragezeichen von der Suchzeichenkette getrennt. [RFC94]

Im Bereich des $\langle searchpart \rangle$ 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]

8.3.2.2 Erzeugung der Such-URLs

Dieser Absatz beschreibt die Erstellung der Such-URLs für Google mit dem Wissen aus Kapitel 8.3.2.1.

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 variierenden 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, wird ein weiterer URL erstellt. Mit diesem 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“. Dabei ersetzt die ASCII-Zeichenkette “%3A“ den Doppelpunkt. [W3S]

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

8.3.2.3 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 “Michael Bazzell“.

Wenn diese Suche ohne Anführungszeichen durchgeführt wird, werden zusätzlich Webseiten vorgeschlagen die den Namen “Mike Bazzell“ anstatt “Michael Bazzell“ beinhalten. Diese erweiterte Suche kann dazu führen, dass unzählige Webseiten vorgeschlagen werden, die nicht unmittelbar etwas 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, Webseiten mit Verweisen auf die Namen “Marco Mustermann“ und “Max Lang“ vorgeschlagen werden. 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. Der 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.

8.3.3 Auswahl der Bibliothek für Serveranfragen

8.3.3.1 Ziele und Anforderungen

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. Dazu müssen Webseiten die JavaScript enthalten auslesbar sein.

8.3.3.2 Lösungsideen

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.

8.3.3.3 Bewertung der Lösungsideen anhand den Anforderungen

Für komplizierte Anfragen an einen Server eignet sich die Request-Bibliothek von Python sehr gut. Darüber hinaus ist der Umgang mit Cookies, Header, et cetera 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 dieselbe Seite keine neue Anmeldung vollzogen werden.

8.3.4 Erstellung des Web Crawlers

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.

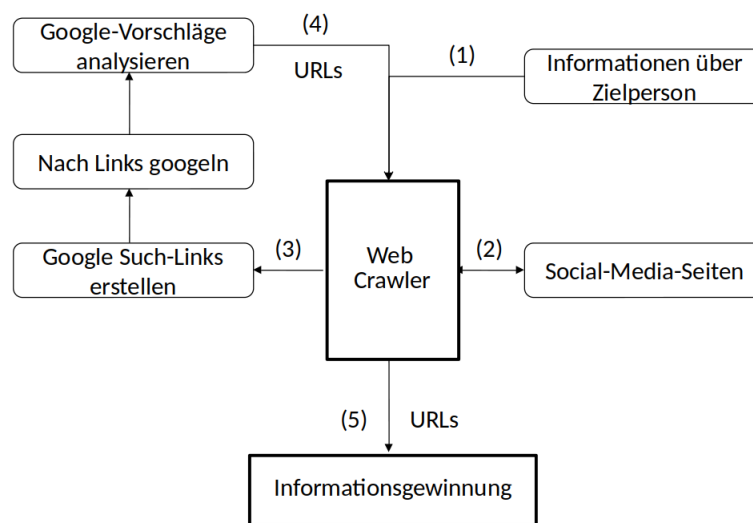


Bild 8.1: Aufbau des entwickelnden Web Crawlers

8.3.4.1 Umgang mit den 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 Behandlung von Social-Media-Seiten zuständig. Damit vollständige Profile angezeigt werden können, loggt sich dieser automatisch auf den entsprechenden Plattformen ein. Die gewonnenen Informationen werden dem Profil der gesuchten Person hinzugefügt.

Es besteht die Möglichkeit, dass nur der Benutzername einer Social-Media-Plattform über die Ziel Person bekannt ist. Weiter Attribute wie Vorname, Nachname und Wohnort sind demnach nicht angegeben. In diesem Fall kann mit einem einzigartigen Benutzernamen nach diesen Attributen auf der entsprechenden Webseite gesucht werden. Einen einzigartigen Benutzernamen wird von den Plattformen Instagram und Twitter verwendet. Aus diesem Grund kann die erweiterte Suche nur auf diesen beiden Plattformen umgesetzt werden.

8.3.4.1.1 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, sobald der Vorname, Nachname und Wohnort eingegeben wurde.

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. Infolgedessen unterscheidet sich die Auswahl der `<input>`-Tags auf den angeforderten Seiten.

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 werden die Felder stattdessen mit Hilfe des Attributes `“class“` gesucht. 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 Attribut `“id“`. Bei Twitter ist das das Attribut `“class[0]“` und bei Xing `“name“`.

8.3.4.1.2 Instagram

Instagram und Twitter verwenden beide einen einzigartigen Benutzername. Dies ist ein bedeutender Vorteil für die Identifikation einer Person. Auf diesen Webseiten kann dadurch eine Person mit ausschließlich einem Benutzername identifiziert werden. Es ist ebenfalls möglich eine Person mit ihrem offiziellen Namen zu suchen und zu finden.

Wenn 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 einer Übereinstimmung wird die Suche nach Kontaktinformationen abgebrochen. Die gefundene Person, sowie die übereinstimmende Daten, werden gespeichert und können später zur Generierung der Phishing-Mail verwendet werden. Eine Übereinstimmung kann beispielsweise dieselbe

Universität sein.

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 Freund 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.

Eine Profilseite mit bekanntem Username wird mit dem folgenden Link angefordert:

<https://www.instagram.com/username/>

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.

8.3.4.1.3 Twitter

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

8.3.4.1.4 Facebook

Facebook bietet ein großes Potential um OSINT zu betreiben. Allerdings hat Facebook optimale Algorithmen zur Erkennung von automatisierten Crawlern entwickelt. Aus diesem Grund wurde das gefälschte 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 verlangt Facebook eine Kopie des Ausweises, ein Bild mit erkennbarem Gesicht und eine Handynummer.

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

8.3.4.1.5 LinkedIn und XING

LinkedIn und XING bieten eine optimale Informationsquelle bezüglich der schulischen und beruflichen Tätigkeit 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. Ein Profil wird untersucht, wenn nur eine Person vorgeschlagen wird. Bei einer Mehrzahl von gefunden 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>

8.3.4.2 Analyse der Google-Suchergebnisse

Zur Analyse der Webseite mit den Suchergebnissen von Google wird der Seitenquelltext benötigt. Mit Hilfe des Quelltextes, können die entsprechenden Links erkannt werden. Der Seitenquelltext 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. Im Ausschnitt des Seitenquelltextes 8.1 ist zu sehen, dass der <div>-Container mit der Klasse "g" einen Hyperlink enthält, was an dem HTML-Tag <a> zu erkennen ist. Dieser Link wird für das Web Scraping benötigt. Deswegen wird genau nach diesem Link gesucht.



Bild 8.2: Container eines Google-Suchergebnisses [LLC19]

Da der `<div>`-Container bei jedem Suchergebnis identisch ist, kann bei jedem Ergebnis nach dem entsprechenden 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. Diese Hyperlinks könnten dazu verwendet werden, um weitere URLs zu finden. Allerdings wird in dieser Arbeit nur die erste Such-Ergebnisseite von Google als Quelle verwendet, damit die Suche analysierbar und bewertbar bleibt.

```
<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]

8.4 Die Personenidentifizierung

8.4.1 Ziele und Anforderungen

Die zu entwickelnden Methoden haben das Ziel die Wahrscheinlichkeit, dass es sich um die gesuchte Person handelt, zu erhöhen. Dabei soll die Laufzeit und Wirksamkeit möglichst optimal sein. Allerdings dürfen hierfür nicht zu viele Informationsquellen ignoriert werden.

8.4.2 Lösungsideen

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 8.2.2 möglich. Dabei wird abhängig von der eingegebenen Information die Suche variiert. Des Weiteren kann durch eine Optimierung des Such-URLs 8.3.2.3 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.

8.4.2.1 Generierung eines Identifikationsschlüssel

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. Das bedeutet, dass es mehrere Personen mit demselben vollständigen Namen geben kann. Um eine Person zu identifizieren, zu welcher keine einzigartigen Informationen bekannt sind, können Kombinationen aus den angegebenen Daten erstellt werden. Diese Kombinationen dienen in diesem Fall als Identifikationsschlüssel. Im folgenden sind alle möglichen Kombinationen aufgelistet.

Vorname, Nachname, Wohnort/Standort;

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.

8.4.2.2 Kontaktanalyse

In den Suchergebnissen kann eine Profilseite von einer Social-Media-Plattformen vorgeschlagen werden. Dabei ist nicht bekannt, ob dieser Benutzer mit der gesuchten Person übereinstimmt. Aus diesem Grund besteht die Möglichkeit einer veränderten Form der Kontaktanalyse, wie in Kapitel 8.10, beschrieben wurde. Hierbei werden die Kontakte ebenfalls durchsucht. Bei einer Übereinstimmung von Daten des Kontaktes mit Attributen der gesuchten Person wird angenommen, dass es sich bei dem gefundenem Profil um die korrekte Person handelt. Die Annahme wird getroffen, da Kontakte bzw. Freunde dadurch eine Beziehung zum Opfer aufweisen. Ein Beispiel hierfür wäre der selbe Standort. Bei keiner Übereinstimmung von Informationen wird die vorgeschlagene Webseite verworfen.

8.4.3 Bewertung der Lösungsideen anhand den Anforderungen

Beide Methoden zur Identifizierung einer Person können eine Verbesserungen der Ergebnisse mit sich bringen. 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 beider Methoden 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.

Das Ergebnis der Kontaktanalyse ist nicht optimal. Es kann nicht davon ausgegangen werden, dass es sich bei einer Übereinstimmung von einem Attribut unmittelbar um die Zielperson handelt. Beim Betrachten der Laufzeit, kann davon ausgegangen werden, dass die Kontaktanalyse deutlich mehr Zeit und Ressourcen benötigt.

Es wird die Methode zur Verwendung eines Identifikationsschlüssels umgesetzt. Dadurch werden ausschließlich die Webseiten behandelt, welche den Schlüssel zur Identifizierung beinhalten. Das bedeutet, dass nur identifizierbare Profile durchsucht werden. Das hat zur Folge, dass manche Profile nicht beachtet werden. Dennoch werden dadurch weitestgehend fehlerhafte Informationen in dem Opferprofil vermieden.

8.5 Umsetzung der Methode zur Verwendung eines Identifikationsschlüssels

Zu Beginn der vorläufigen Inhaltskontrolle werden die Eingaben abgefragt. Dadurch wird erkannt, zu welche Daten Informationen vom Benutzer eingegeben wurden. Anschließend werden mit diesen Daten alle möglichen Kombinationen aus Kapitel 8.4.2.1 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, Nachnamen und beispielsweise Wohnort in

dem Webseitentext befindet. Wenn diese Abfrage korrekt ist, wird die Webseite weiter behandelt und es kann nach Information gesucht werden.

8.6 Herausfiltern von wichtigen Informationen auf einer Webseite

8.6.1 Ziele und Anforderung

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.

8.6.2 Lösungsideen

Zu Beginn werden Wortsammlungen erstellt. Diese Wortsammlungen sind Listen, welche aussagekräftige Schlüsselwörter enthalten und nach Themen kategorisiert sind. Beispiele für den Inhalt dieser Listen sind alle Hochschulen- und Universitätsnamen in Deutschland, Berufsbezeichnungen und Tätigkeiten, Hobbybezeichnungen sowie alle Städte und Gemeinden in Deutschland. Die Wortsammlungen werden mit Hilfe von bekannten Listen im Internet eigenständig befüllt. Als Informationsquelle dienen alle öffentlich frei zugänglichen Quellen, die hilfreiche Informationen enthalten. Die erzeugten Listen werden in den kommenden Schritten verwendet, um wichtige Informationen herausfiltern zu können.

Die Grundidee zum Erkennen von wichtigen Informationen 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 und in einer Liste gespeichert. Diese Liste kann darauf mit den erzeugten Wortsammlungen verglichen werden. Bei einer Übereinstimmung

eines Schlüsselwortes wird das Wort mit der entsprechenden Kategorie vorgemerkt und später in die verwendete Speicherstruktur eingetragen.

Eine weitere Methode zur Textanalyse ist der Vergleich von Zeichenketten. Dabei wird der vorliegende Webseitentext in einen String umgewandelt. Anschließend kann kontrolliert werden, ob Elemente aus den Wortsammlungen in diesem Text vorkommen. Auch hier kann bei einer Übereinstimmung das Wort mit entsprechenden Kategorie gespeichert werden.

8.6.3 Bewertung der Lösungsideen anhand den Anforderungen

Für die Methode, bei der ausschließlich Zeichenketten verglichen werden, besteht die Möglichkeit, alle Elemente der Wortsammlung zu erkennen. Dabei ist nicht relevant, aus wie vielen Wörtern und Zeichen ein Element besteht. Allerdings kann diese Methode zu falschen Ergebnissen führen. Zum Beispiel bei der Suche nach dem Wort "test". Hierbei wird eine Übereinstimmung mit dem Wort "testament" gefunden, da die Zeichenfolge des gesuchte Elements in diesem Wort vorkommt. Trotzdem haben beide Elemente eine unterschiedliche Bedeutung.

Ein Nachteil der Schlüsselwortgenerierung ist die festgelegte Anzahl an Wörter, aus dem ein Schlüsselwort besteht. Damit Elemente mit mehreren Wörtern gefunden werden können, müssten N-Gramme erstellt und den Schlüsselwörtern hinzugefügt werden. Andernfalls werden ausschließlich Wörter, bestehend aus einem Wort, in Betracht gezogen. Im Gegensatz zur vorherigen Methode, kann mit dem Verfahren zur Schlüsselwortgenerierung die Anzahl des Vorkommens eines Wortes gezählt werden. Diese Anzahl ist wichtig für die Auswahl der Information. Darüber hinaus entstehen hierbei keine Fehler durch falsche Interpretationen.

Aus den erwähnten Gründen wird für das Herausfiltern von wichtigen Informationen die Methode der Schlüsselwortgenerierung verwendet.

8.7 Umsetzung der Methoden zum Herausfiltern von wichtigen Informationen

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 zuerst 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 zwei 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 Sammlung mit Schlüsselworten soll manuell erstellbar und beliebig erweiterbar sein. Die Anwendung muss ohne aufwendige Zugriffe von diesen Listen lesen können. Als mögliches Speichervariante dient eine CSV-Datei oder eine SQL-Datenbank.

Die SQL-Datenbank ist ein komplexeres System. Infolgedessen werden aufwendige Zugriffe 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.

Die Informationsgewinnung ist abhängig von diesen Wortsammlungen. Das bedeutet, dass nur die Informationen herausgefunden werden können, welche als Element in einer Wortsammlung vorkommt. Infolgedessen hat die Umsetzung der Wortsammlungen einen hohen Stellenwert.

8.7.2.1 Umsetzung der Wortsammlungen

Für eine sinnvolle Informationsgewinnung werden die Wortsammlungen kategorisiert. Dabei entsprechen die Kategorien einem Teil der zu suchenden Personenattributen, wie in Bild 8.8 dargestellt wird. Demzufolge gibt es die Kategorie “Tätigkeiten“, “Hobbys“, “Institutionen“ sowie “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. Befüllt werden die Sammlungen aus öffentlich zugänglichen Listen im Internet. Die Tabelle 8.7.2.1 zeigt die Kategorien der Wortsammlungen und die dazugehörigen Informationsquellen. Zur Veranschaulichung ist ein Ausschnitt der Liste mit den Tätigkeiten im Anhang A beigelegt.

Kategorie	Informationsquelle
Tätigkeiten	http://planet-beruf.de/schuelerinnen/mein-beruf/berufe-von-a-z
Hobbys	https://hobbyfuchs.org/hobbyliste-a-z ; https://hobbeasy.de/hobbys-finden-liste ; https://www.langeweile-tipps.de/hobby-finden-die-ultimate-liste-der-hobbys
Institutionen	https://de.wikipedia.org/wiki/Liste_der_Hochschulen_in_Deutschland ; https://www.firmenfinden.de (Firmennamen nur von Kreis Bodensee und Ravensburg)
Städte und Gemeinden	https://www.deutschland-auf-einen-blick.de/index.html

Tabelle 8.1: Informationsquellen der Wortsammlungen

8.7.3 Automatisierte Schlüsselwortgenerierung

8.7.3.1 Ziele und Anforderungen

Aus einem vorliegenden Webseitentext sollen Schlüsselwörter generiert werden. Dafür muss jedes Wort in diesem Text betrachtet und als Schlüsselwort erkannt werden. Des Weiteren

soll die Möglichkeit bestehen, dass die Anzahl, aus wie vielen Wörtern ein Schlüsselwort besteht, festgelegt werden kann.

8.7.3.2 Lösungsideen

Möglichkeiten zur automatisierten Schlüsselwortgenerierung stellen die Verfahren Rapid Automatic Keyword Extraction (RAKE) und Automatic Keyword Extraction mit NLP dar.

8.7.3.2.1 RAKE

RAKE ist eine sehr effiziente Methode zur Schlüsselwortgenerierung. 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 aufgeteilt. Die entstehenden Fragmente, welche als mögliche Schlüsselwörter dienen, werden in ein Array gespeichert. 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 berechnet. 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 das gemeinsame Auftreten mit sich selbst und anderen Schlüsselwörtern. In der Tabelle 8.7.3.2.1 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 8.7.3.2.1 ablesen. Allerdings muss hier die Reihe und Spalte des jeweiligen Wortes abgeglichen 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 dargestellten 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 8.2: 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

8.7.3.2.2 Automatic Keyword Extraction mit NLP

Bei dieser Methode wird ein 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 “Bad Waldsee“ herauszulesen.

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

8.7.3.3 Bewertung der Lösungsideen anhand den Anforderungen

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 Webseitentext 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.

Die Methode zur automatisierten Schlüsselwortgenerierung mit NLP bringt dagegen eine 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.

8.7.4 Umsetzung des 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 dieser Liste werden die “stopwords“ der deutschen als auch der englischen Sprache

herausgefiltert. Dadurch verringert sich die Anzahl der zu suchenden Wörter. Anschließend werden Bigramme aus dem zerlegten Text erstellt und der Liste mit den einzelnen Wörtern hinzugefügt. Dadurch können Elemente bestehend aus einem Wort sowie aus zwei Wörtern gefunden werden.

Im nächsten Schritt wird die Liste mit den entsprechenden Wortsammlungen verglichen. Dabei werden die übereinstimmenden Wörter dem korrekten Personenattribut hinzugefügt. Die Wortsammlung der Institutionen wird nicht mit der Schlüsselwortliste verglichen. Die Problematik besteht darin, dass die Institutionsnamen aus verschiedenen vielen Wörtern bestehen können. Für die Erkennung solcher Namen müsste eine große Anzahl von N-Grammen erstellt und der Schlüsselwortliste hinzugefügt werden. Das würde allerdings zu einem großen Laufzeitnachteil führen. Demnach wird der Institutionsname in dem unformatierten Webseitentext gesucht. Der Text wird lediglich auf das Vorkommen der korrekten Zeichenkette kontrolliert. Dadurch kann das Vorkommen der Institution auf einer Webseite nicht gezählt werden. Es kann nur die Anzahl der Webseiten, welche diese Zeichenkette beinhalten, festgestellt werden. Dennoch ist die Fehlerrate gering, da die langen Zeichenketten oft einzigartig sind.

8.7.5 Herausfiltern von Geburtsjahren

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 8.7.4 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. Die Zahl wurde dabei auf fünfzehn festgelegt, da die Mehrzahl der geschriebenen Sätze eine Anzahl von 12-15 Wörtern enthält. [Sei69] Falls dabei das Wort "Geburtsdatum", "Alter", "geboren", "Geburtsort", "born", oder "birth" vorkommt, wird das entsprechende Jahr als Geburtsjahr der Zielperson festgelegt. Die Suche nach den Signalwörtern wird mit Hilfe des folgenden regulären Ausdrucks durchgeführt: `r"((geburtsdatum)|(alter)|(geboren)|(geburtsort)|(born)|(birth))"`.

Das Bild 8.5 zeigt die Funktion des Algorithmus zur Suche nach dem Geburtsjahr. Dabei

wurde das Jahr “1995“ gefunden. Somit werden 15 Wörter vor und nach dieser Jahreszahl kontrolliert. Falls der Algorithmus vorzeitig eine Übereinstimmung mit dem regulären Ausdruck findet, wird die Suche, wie in diesem Beispiel, abgebrochen. Der rote Pfad im Bild 8.5 zeigt auf das korrekt Signalwort “geboren“. Dadurch wird die Suche nach dem fünften Schritt erfolgreich beendet und die Jahreszahl “1995“ dem Opferprofil als Geburtsjahr hinzugefügt.

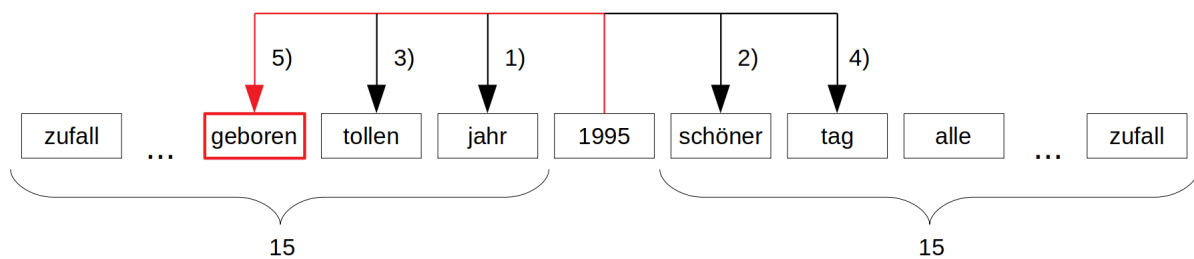


Bild 8.5: Funktion des Algorithmus zur Suche nach dem Geburtsjahr

8.7.6 Herausfiltern von E-Mail-Adressen

Zu Beginn wird der unformatierte Webseitentext in Textfragmente zerlegt. Das Leerzeichen dient dabei als Trennelement. Anschließend werden die erzeugten Fragmente mit dem regulären Ausdruck, $r = "(.*((@)|(at)).*(de|com|net)).*"$, nach einer gültigen E-Mail-Adresse durchsucht. Bei einer Übereinstimmung des regulären Ausdrucks wird der korrekte Teilstring und somit die E-Mail-Adresse ausgelesen.

8.7.6.1 Auswahl 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. Diese

Grenze wurde aus den Ergebnissen von zahlreichen selbst durchgeführten Tests auf die Zahl 0,5 % festgelegt. Dafür wurden Wertungen für verschiedene E-Mail-Adressen mit dem dazugehörigen Namen berechnet. 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 “mustermax“ und “martafrau“. Anschließend werden die lokalen Namen der E-Mail-Adressen mit dem erzeugten Namensstring der Zielperson verglichen. Dabei erreicht die lokale Namen *mustermax* eine prozentuale Übereinstimmung von 0,73 % mit dem Namensstring und *martafrau* 0,27 %. Da die Prozent-Grenze bei 0,5 % beträgt, wird die zweite E-Mail-Adresse verworfen.

8.8 Auswahl der gefundenen Information

In diesem Abschnitt werden die Methoden zur Auswahl von den gefundenen Informationen beschrieben. Das hat den Grund, dass alle gewonnenen Schlüsselwörter einer Webseite in einer Liste gespeichert sind. Jedoch kann nur eines dieser Schlüsselwörter für die Generierung der Phishing-Mail verwendet werden.

8.8.1 Ziele und Anforderungen

Ziel ist es ein Schlüsselwort aus jeweils einer Liste der gefundenen Kategorien zu bestimmen. Dabei soll beachtet werden, wie oft eine Element auf einer Webseite und auf unterschiedlichen Webseiten vorkommt.

8.8.2 Lösungsidee

Es kann eine Methode entwickelt werden, welche einen Score für die Bewertung eines Elementes berechnet. Hierfür wird eine prozentuale Wertung für das Vorkommen eines Wortes in einer Liste mit der Formel 8.1 berechnet, nachdem eine Seite vollständig durchsucht wurde.

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

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“ siebenmal öfter als das Wort “Fechten“ auf der Webseite vorkommt. Solch eine Liste wird für jede durchsuchte Seite erstellt. Nachdem alle Webseiten durchsucht wurden, werden alle erstellten Listen zu einer zusammengefügt. Dabei bleibt die Struktur bestehen, damit erkannt wird, welche Wörter auf unterschiedlichen Webseiten vorkommen. 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 eine Matrix mit allen Wertungen erstellt. Eine Matrix für das hier verwendete Beispiel ist in Tabelle 8.8.2 aufgezeigt. Die Spalten enthalten alle Elemente einmal, die in der Liste vorkommen können. Die Zeilen entsprechen der Anzahl der durchsuchten Webseiten.

Mit der nachfolgenden Formel 8.2 wird für jedes Element ein endgültiger Score aus der Matrix berechnet. Dafür werden die Spalten summiert und durch die Anzahl der durchsuchten Webseiten geteilt.

	Fußball	Basketball	Fechten	Volleyball
Webseite 1	0.7	0.2	0.1	0
Webseite 2	0.5	0	0	0.5

Tabelle 8.3: Beispiel-Matrix

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

mit

m = Anzahl aller möglichen Elemente (Spalten) (8.2)

n = Anzahl der durchsuchten Webseiten (Zeilen)

i = Spalte

j = Zeile

In dem aufgezeigten Beispiel würde das zu folgendem Ergebnis 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 es das Wort "Fußball". Falls zwei Wertungen gleich hoch sind, wird das Wort, welches als erstes in der entsprechenden Liste vorkommt, ausgewählt.

8.8.3 Bewertung der Lösungsidee anhand den Anforderungen

Die vorgestellte Methode zur Berechnung eines Scores, bringt alle Anforderungen mit sich. Es wird nicht nur das Vorkommen eines Elements in einer Liste gezählt, sondern berücksichtigt wie oft es auf einer Webseite und auf wie vielen Webseiten es vorkommt. Dadurch können Fehler bei der Gewichtung eines Elements vermieden werden. Das hat den Grund, dass ein Schlüsselwort auf einer Webseiten öfter auftreten kann. Im Fall dass ein Element auf einer Webseite sehr oft gefunden wird, könnte mit einer Methode

bei der ausschließlich das Vorkommen gezählt wird ein Fehler entstehen. Somit würde ein häufig auftretendes Schlüsselwörter auf einer Webseite andere überstimmen, obwohl diese ebenfalls oft und auf verschiedenen Webseiten zu finden sind. Infolgedessen wird die Methode zur Berechnung eines Scores umgesetzt.

8.9 Umsetzung der Methode zur Auswahl der gefundenen Information

Zu Beginn wird unterschieden ob ein Score für die Kategorie „Institution“ oder für eine anderer berechnet werden soll. Das hat den Grund, dass Liste bei diesen Kategorien unterschiedlich sortiert werden. Im Fall dass der Score für die Institutionen ermittelt werden soll, wird die Liste mit den gefundenen Elementen der Länge nach sortiert. Somit ist das Element mit der längsten Zeichenkette an erster Stelle. Das hat den Grund, dass bei der Suche nach Institutionen lediglich die Zeichenketten gesucht werden. Somit kann es sein, dass ein langer Institutions-Name unbeabsichtigt einen kürzen beinhaltet. Ein Beispiel hierfür ist die Firma “MF Musterfirma GmbH & Co. KG“, welche zusätzlich einen anderen Firmennamen wie “Musterfirma GmbH“ einschließt. Durch die eben erwähnte Sortierung wird dieser Fehler umgangen.

Falls es sich allerdings um eine andere Kategorie handelt, werden die Liste nach der Häufigkeit wie oft ein Element auf einer Webseite vorkommt sortiert. Somit wird bei einem gleichen Score das Element gewählt, welches ein höheres Vorkommen auf einer Webseite hat.

Mit der Programmbibliothek NumPy wird eine Matrix erstellt und mit Nullen befüllt. Anschließend werden die errechneten Scores an den entsprechenden Positionen eingesetzt. Daraufhin werden die Spalten zusammengerechnet und verglichen. Das Element, welches der Spalte mit dem höchsten Wert entspricht, wird ausgewählt. Im Listing 8.2 wird der Algorithmus zur Berechnung des finalen Scores und der Auswahl des entsprechenden Elements dargestellt.

```
score = 0
element_with_highest_score = ""
#every column
for k in range(0, numpy.size(matrix, 1)):
    #current_score = sum of elements of a column
    current_score = 0
    #every row
    for i in range(0, numpy.size(matrix, 0)):
        current_score = current_score + matrix[i][k]
    current_score = current_score / len(list3)
    if current_score > score:
        score = current_score
        element_with_highest_score = instances[k]
```

Listing 8.2: Algorithmus zur Berechnung des Scores und der Auswahl des Elementes mit dem höchsten Score

8.10 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. Durch die erwähnte Methode können weitere Informationen gewonnen werden.

8.10.1 Ziele und Anforderungen

Ergänzend zu der Personensuche sollen Kontaktinformationen gefunden werden. Als Informationsquelle kann jede Social-Media-Plattformen dienen, welche es ermöglicht die Kontakte der gesuchten Person anzuzeigen. Des Weiteren sollen Informationsquellen ohne konfigurierbare Sicherheitsvorkehrungen bevorzugt werden, da ansonsten mögliche Kontakte nicht angezeigt werden können. Um die Funktion der Kontaktanalyse aufzuzeigen ist es ausreichend diese auf einer Social-Media-Webseite durchzuführen.

Eine gefunden Kontaktinformation muss äquivalent zu einem Attribut des Opfers sein.

Somit kann davon ausgegangen werden, dass beide Personen die selben Interessen in dem Bereich der gefundenen Information haben.

8.10.2 Lösungsideen

Als mögliche Informationsquellen für Kontakte zählen die Social-Media-Seiten von Facebook, LinkedIn, Instagram, Twitter und XING.

8.10.3 Bewertung der Lösungsideen anhand den Anforderungen

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 Benutzern nicht erkennbar. Facebook, Twitter und Instagram bieten die Möglichkeit, die Kontakte der gesuchten Person anzuzeigen. Allerdings wird dafür ein Account benötigt. Somit eignen sich die Seiten Twitter, Xing, Facebook und Instagram. Wie in Kapitel 8.3.4.1 beschrieben, wird kein Facebook-Account angelegt. Dadurch ist es nicht möglich, Kontakte auf dieser Webseite anzuzeigen. Um die Funktion der Methode aufzuzeigen, wird ausschließlich die Webseite Instagram verwendet, da sie alle Anforderungen mit sich bringt und im Ranking der größten sozialen Netzwerke an aktiven Nutzern vor Twitter steht. [Dat19]

8.11 Umsetzung der Instagram-Kontaktanalyse

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, welche 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. Ein Beispiel hierfür kann die selbe Universität oder der selbe Wohnort sein. Sobald dies gefunden wurde, kann die Suche beendet werden. Wenn keine Übereinstimmung der Profile gefunden wurde, wird dem erstellten Opferprofil keine Kontaktformation hinzugefügt.

8.11.1 Auslesen der Kontakte

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 “Dieses 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 8.6 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.

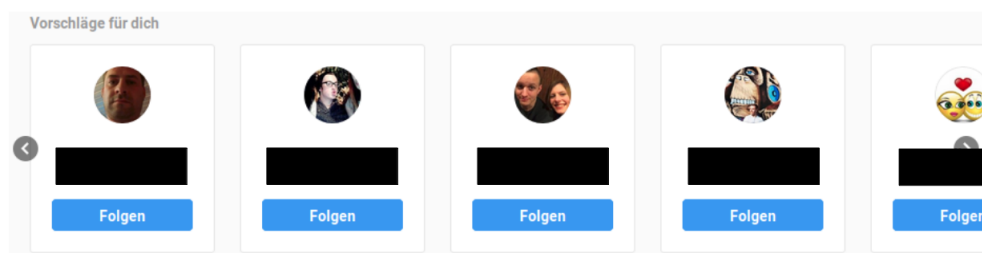


Bild 8.6: Container mit Profil-Vorschlägen [INS19]

Falls es sich um ein öffentlich frei zugängliches Profil handelt, kann eine Liste der abonnierten Kontakte angezeigt werden. Diese sind über ein scrollbares Pop-Up-Fenster, wie in Bild 8.7 dargestellt, einsehbar. Vergleichbar zur Methode bei einer privaten Profilseite

wird hier ebenfalls Schritt-für-Schritt durchgescrollt. Dadurch wird jeder Kontakt geladen. Somit können alle Links, welche zu den entsprechenden Profilseiten führen, ausgelesen werden.

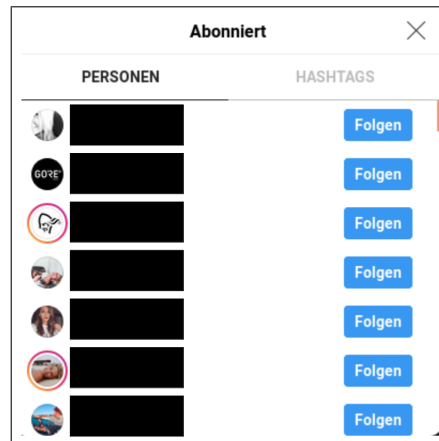


Bild 8.7: Pop-up-Fenster mit abonnierten Profilen [INS19]

Die Herausforderung besteht darin, dass nicht zu schnell gescrollt werden darf. Andernfalls werden keine weiteren Kontakte geladen. Aus diesem Grund wird ein Algorithmus 8.3 verwendet, welcher einem menschlichen Verhalten ähneln soll und somit langsam und schrittweise nach unten scrollt. Hierfür 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 Scroll-Vorgang wird das Fenster nur ein sechstel des möglichen Bereichs nach unten geblättert. Dadurch werden weitere Profile geladen.

In den nächsten Schritten wird das Fenster jeweils ganz nach unten verschoben. Dadurch werden alle Profile geladen. Sobald alle internen 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 gefundene 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_class_name("isgrP")
# 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 8.3: Algorithmus zum Herunterscrollen des Pop-up Fensters

8.12 Speicherung der gewonnenen Daten

8.12.1 Ziele und Anforderungen

Die gespeicherten Daten werden von verschiedenen Klassen benötigt. Aus diesem Grund muss es möglich sein, dass andere Klassen auf die Speicherstruktur zugreifen können. Zusätzlich wird eine gute Struktur vorausgesetzt, damit einzelne Attribute der Person ausgewählt werden können. Des Weiteren muss die Speichervariante in Python implementierbar sein und unnötige Lese- und Schreibzugriffe sollen vermieden werden. Es ist nicht notwendig, dass die Daten nach Programmende abrufbar sind. Infolgedessen wird keine externe Speicherung in einer Datenbank oder in einer Datei vorausgesetzt.

8.12.2 Lösungsideen

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.

8.12.3 Bewertung der Lösungsideen anhand den Anforderungen

Die Umsetzung jeder einzelnen Variante ist mit Python möglich. Für die Verwendung einer SQL-Datenbank spricht die gute Speicherstruktur. Allerdings sind mit einer solchen Datenbank aufwendigere Speicher- und Lesevorgänge verbunden. Die externe Speicherung in einer Datei wie CSV oder TXT ist keine Anforderung. Darüber hinaus müssten diese Dateien verschlüsselt werden, damit sie vor fremden Zugriffen geschützt sind. Dennoch werden die gewonnenen Daten in einem Personenobjekt gespeichert. Unnötige Speicher- und Lesezugriffe fallen dadurch weg. Darüber hinaus lässt sich das Personenobjekt einfach an die entsprechenden Klassen übergeben.

8.13 Umsetzung der Personenklasse

Die gewonnenen Daten werden in einem Personenobjekt, wie in Bild 8.8 dargestellt, gespeichert. Dabei werden die vom Anwender eingegebenen Daten direkt in das Personenobjekt übertragen. Zusätzlich werden die gefundenen Personenattribute hinzugefügt. Diese sind zu Beginn in Form einer Liste gespeichert. Sobald das häufigste Element mit der Methode 8.9 ermittelt wurde, wird nur das entsprechende Element gespeichert.

Zu einem angegebenen Attribut kann eine zusätzliche Information gefunden werden. So ist es denkbar, dass beispielsweise zu einem bekannten Wohnort ein weiterer Ort gefunden wird. Dies ist für die Generierung der Phishing-Mail wichtig. Die Zielperson hat möglicherweise einen höheren Bezug zu dem gefundenen Ort. Aus diesem Grund wird bei der Informationsauswahl für die Phishing-Mail die gefundene Information den angegebenen Daten bevorzugt.

Falls bei der Kontaktanalyse ein übereinstimmendes Profil gefunden wurde, kann diese

Information in dem Attribute “Gefundene Kontaktinformation“ gespeichert werden. Hierbei wird zuerst der vollständige Kontaktname und anschließend die übereinstimmende Information gespeichert. Ein Beispiel hierfür ist ['Max Mustermann', 'Fußball'].

PERSON
Vorname Nachname Geschlecht Wohnort/Standort Geburtsjahr Institution Instagram-Benutzername Facebook-Benutzername Twitter-Benutzername E-Mail-Adresse Gefundene Tätigkeit Gefundenes Hobby Gefundene Institution Gefunden E-Mail-Adresse Gefundene Orte Gefundene Kontaktinformation Durchsuchte Links

Bild 8.8: Personenklasse

9 Implementierung der Phishing-E-Mail

In diesem Kapitel befindet sich die Umsetzung zur Erstellung einer Phishing-Mail. Dabei wird zuerst auf die Implementierung des Algorithmus zur Adressgenerierung eingegangen. Anschließend werden die E-Mail-Muster vorgestellt. Im letzten Abschnitt wird die Umsetzung zum Versenden einer Phishing-E-Mail aufgezeigt.

9.1 Umsetzung des Algorithmus zur Generierung von 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 generiert werden kann. Jedoch wird in der Anwendung keine Phishing-Mail an eine dieser Adressen, welche mit der Methode erzeugt wurden, versendet. Es dient ausschließlich als Beweis für die Durchführbarkeit dieser Methode.

9.1.1 Funktion des 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 bürgerliche Namen verwendet. [Med17] Aus diesem Grund verwendet der Algorithmus die Personattribute Vorname, Nachname und zusätzlich 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 des Geburtsjahrs;

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 von Anordnungen der bekannten Daten generiert. Dabei wird zusätzlich auf die Reihenfolge der Anordnungen geachtet. Somit ist das Element “Max Mustermann” und “Mustermann Max” eine eigene Kombination. Es werden zusätzlich die selben Kombinationen mit den Trennzeichen “.”, “_” und “-” erstellt und in einer Liste gespeichert. Das Ergebnis einer Permutation über die Attribute “Vorname” und “Nachname” werden nachstehen dargestellt.

[Vorname], [Nachname], [Vorname & Nachname], Vorname & '.' & Nachname], [Vorname & '_' & Nachname], [Vorname & '-' & Nachname], [Nachname & Vorname], [Nachname & '.' & Vorname], [Nachname & '_' & Vorname], [Nachname & '-' & Vorname]

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 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

Tabelle 9.2: Erzeugte E-Mail-Adressen zu den Attributen “Marco“, “Lang“ und “1995“

9.2 Umsetzung 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 und Gefühle wie Ängste und Freude auslösen. Aus diesem Grund muss für jede Kombination der gewonnenen Daten ein Muster zur Verfügung stehen. Zusätzlich stellt sich die Frage, wie die E-Mail-Texte möglichst passend kategorisiert werden können.

9.2.1 Kategorisierung der E-Mail-Muster

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 Umgangssprache und ein berufliches in einer formelleren Sprache erstellt wird. Diese beiden Kategorien haben weitere Unterkategorien, welche verschiedene Kombinationen aus den personenbezogenen Daten verwenden. Bei beiden Kategorien werden menschliche Gefühle und Verhaltensmuster ausgenutzt.

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. Andernfalls wird ein privates Muster verwendet.

9.2.1.1 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 die kursiv geschriebenen Wörter in den Bildern mit den gewonnenen Daten über die Zielperson ersetzt.

Das Bild 9.1 beschreibt ein berufliches Muster, welches die Personenattribute Nachname

und Institution verwendet. Dabei wird durch diesen Inhalt ein Pflichtgefühl bei der Zielperson ausgelöst, wodurch die Wahrscheinlichkeit steigt, dass sie auf den entsprechenden Link klickt. Im Bild 9.2 ist ein Muster mit den Attributen Nachname und Tätigkeit aufgezeigt. Hierbei wird die Neugier des Opfers angesprochen. Des Weiteren werden positive Gefühle bei der Zielperson ausgelöst, da die Chance auf eine neue Herausforderung besteht. Das Bild 9.3 zeigt das Muster mit denselben 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. Dadurch können bekannte Verhaltensmuster und Abläufe im Alltag missbraucht werden. Im letzten Bild wird das berufliche Muster für die Daten Nachname, Tätigkeit und Institution beschrieben, wobei ebenfalls die Neugier des Opfers angesprochen wird.

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* gesucht - Im Auftrag des BRD

Hallo *Herr Mustermann*,

die Bundesrepublik Deutschland sucht einen kompetenten *Mustertätigkeit*. Haben Sie Interesse an einer neuen Herausforderung unter optimalen Arbeitsbedingungen?

Im Anhang finden Sie die offizielle Stellenausschreibung mit den dazugehörigen Voraussetzungen und Gehaltsstufen.

Ihr Karriere-Team der Bundesrepublik Deutschland

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

SUBJECT: Feedback zur Ausarbeitung

Hallo *Herr Professor Mustermann*,

wie besprochen befindet sich im Anhang meine vorläufige Ausarbeitung. Könnten Sie diese bitte 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 folgendem 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

9.2.1.2 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 9.2.1.1.

Im ersten Bild wird ein privates E-Mail-Muster beschrieben, welche die gewonnenen Kontaktinformationen verwendet. Genaugenommen ist das der Kontaktnamen und die Kontaktinformation. Diese Kontaktinformation ist identisch zu einer Information über die Zielperson. Hierbei wird ein Vertrauen durch einen bekannten Kontaktnamen zum Opfer aufgebaut. Zusätzlich wird das Interesse und die Neugier durch die Verwendung der übereinstimmenden Information sowie die Hilfsbereitschaft dieser Person missbraucht. Das nächste Bild beschreibt ein Muster, bei dem die Personendaten Vorname und Hobby verwendet werden. Dabei wird das persönliche Interesse der Zielperson angesprochen. Das Bild 9.7 zeigt die Verwendung des Vornamens und des Geburtsjahrs in einem E-Mail-Muster. Durch die entstehende Neugier, wissen zu wollen welche Personen im gleichen Jahr geboren sind, wird das Opfer dazu verleitet den Anhang zu öffnen. Das letzte Muster verwendet die Personenattribute Vorname und Ort. Hierfür wird der gefundene Ort verwendet und das Gefühl von Hunger ausgenutzt.

SUBJECT: Fragen bzgl. *Kontaktinformation*

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

Kontaktnamen

Bild 9.5: Ein privates E-Mail-Muster mit den Personenattributen Vorname, Kontaktnamen 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*,
derzeit 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 Streetfood-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 Absender-Adresse benötigt. Sehr große Provider wie Gmail oder Yahoo sind dafür nicht optimal. Das hat den Grund, dass viele Spammer diese Provider in der Vergangenheit verwendet haben. Dadurch werden diese Adressen gerne öfter überprüft. Kleine Provider sind dagegen nicht so bekannt. [Baz18]

9.3.1 Ziele und Anforderungen

Es wird ein Account bei einem kleinen Provider benötigt, damit er großen Social-Media-Plattformen unbekannt ist. Dabei sollen keine Verbindungen zu einem anderen Account oder einer Person bestehen. Zudem soll es ein freemail Account sein.

9.3.2 Lösungsideen

GMX, WEB.DE, 1&1 und Freenet sind mögliche Provider für die Erstellung eines gefälschten Accounts.

9.3.3 Bewertung der Lösungsideen anhand den Anforderungen

Der Provider GMX ist einer der wenigen Anbieter, bei dem 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 verwendeten Profil oder einer Person besteht. Des Weiteren ist dieser kleine Provider großen Plattformen wie Facebook fremd und wird dadurch weniger streng überprüft. Da der Provider von GMX alle Anforderung erfüllt, kann dieser vorzeitig ausgewählt werden. Infolgedessen werden die Phishing-Mails mit einer GMX-Mail-Adresse versendet. [Baz18]

9.4 Umsetzung zum Versenden einer Phishing-E-Mail

Mithilfe 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. Als nächstes wird die Verbindung mit dem GMX-Mailserver aufgebaut. Nach einer erfolgreichen Anmeldung kann eine E-Mail versendet werden.

10 Evaluation der Implementation

10.1 Validierung des Gesamtkonzeptes

Das Gesamtkonzept dieser Anwendung wurde entsprechend den Anforderungen umgesetzt. Eine große Herausforderung ist die Identifizierung einer Person. Dazu wurde die Methode zur Generierung von Identifikationsschlüsseln erstellt. Darüber hinaus wurden die Sucher-URLs optimiert, damit die Suchergebnisse reduziert und verbessert werden. Dennoch besteht die Möglichkeit der Verwechslung einer Person, wenn sich die gefundene Profile sehr ähneln.

Zum Herausfiltern von wichtigen Informationen wurden Schlüsselwörter aus dem Text erzeugt und mit den Elementen aus den Wortsammlungen verglichen. Dabei bestehen alle Elemente, außer die der Wortsammlung Institution, aus ein bis zwei Wörtern. Das bedeutet, es können nur Schlüsselwörter bestehend aus einem Wort oder zwei Wörtern gefunden werden. Dagegen kann eine Institution mehrere Wörter enthalten. Allerdings kann in diesem Fall die Häufigkeit des Vorkommens einer Institution auf einer Webseite nicht erkannt werden. Somit besteht bei beiden Fällen die Gefahr von Fehlinterpretation oder Missachtung einer Information.

Für die Bestimmung, welche Daten verwendet werden, wurde ein Algorithmus entwickelt. Dieser berechnet unter Beachtung von bestimmten Kriterien einen Score. Anschließend wird das Element mit der höchsten Score ausgewählt. Diese Berechnung kann im Fall, dass nur ein Element einer Kategorie auf einer Webseite gefunden wurde, zu Problemen führen. Bei dieser Ausnahme wird das eine Element sehr hoch gewichtet. Dadurch kann es zu einer fehlerhaften Auswahl kommen. Jedoch ist die Berechnung der Wertung unter Berücksichtigung der Kriterien nötig, um einen dauerhaften Auswahlfehler zu vermeiden. Die Phishing-E-Mails werden mit Mustern erzeugt. Dadurch enthält die E-Mail einen

sinnvollen Inhalt mit einer korrekten Grammatik. In einzelnen Fällen kann es jedoch zu sonderbaren Formulierungen kommen.

10.2 Beschreibung und Motivation der Testfälle

Bei den Testfällen wurden verschiedene Personen mit unterschiedlichen Daten gesucht. Dazu wurde für jede Person eine Einverständniserklärung eingeholt. Somit besteht jeder einzelne Testfall aus einer Suche nach einer realen Personen. Dennoch werden Pseudonyme zum Schutz der persönlichen Daten verwendet.

10.2.1 Testfall 1

Im ersten Testfall wurde nach der Person mit dem Namen “Marco Lang“ und dem Wohnort Tettnang gesucht. Dabei ist zusätzlich der Instagram-Benutzername dieser Zielperson bekannt.

10.2.1.1 Ergebnisse

```
Vorname: marco
Nachname: lang
Wohnort/Standort: tett nang
Geburtsjahr: 1995
Ort: tett nang
Tätigkeit: mauerer
Hobby: fitness
Institution: None
E-Mails: []
Kontaktinformation: ['sophie', 'fitness']

Phishing-Mail:
Betreff: Fragen bzgl. Fitness
Hi Marco,
hier ist Sophie. Bezüglich Fitness hätte ich noch ein paar fragen an
dich...
Könntest du zufällig in den Anhang schauen und bewerten was ich da so
rausgesucht habe?

Grüße,
Sophie
```

Bild 10.1: Programmausgabe zum Testfall 1 - “Marco Lang“, “Tett nang“ und “Instagram-Benutzername“

10.2.2 Testfall 2

Für diesen Testfall wird nach der Person “Anita Schmidt“ gesucht. Der Wohnort ist bei dieser Suche keine Stadt, sondern die Gemeinde Heidesheim aus dem Bundesland Rheinland-Pfalz. Es werden keine zusätzlichen Personendaten angegeben.

10.2.2.1 Ergebnisse

```
Vorname: anita
Nachname: schmidt
Wohnort: heidesheim
Geburtsjahr: None
Ort: mainz
Tätigkeit: student
Hobby: motorrad
Institution: hochschule mainz
E-Mails: []
Kontaktinformation: []

Phishing-Mail:
Betreff: Rückmeldung - Hochschule Mainz
Hallo Frau Schmidt,
leider ist und ein Fehler unterlaufen. Aus diesem Grund müssen sie sich
erneut zurückmelden. Um den Vorgang zu beschleunigen, klicken Sie bitte
auf den folgenden Link.
https://badlink.com

Mit freundlichen Grüßen

Ihr Studentenservice der Hochschule Mainz
```

Bild 10.2: Programmausgabe zum Testfall 2 - "Anita Schmidt" und "Heidesheim"

10.2.3 Testfall 3

Für diesen Testfall ist die Zielperson "Klaus Maier". Hierbei wird zweimal nach der gleichen Person gesucht. Allerdings mit zwei unterschiedlichen Orten. Der erste Ort ist Meckenbeuren und entspricht dem Arbeitsort. Im zweiten Fall wird der Wohnort Tettnang verwendet.

10.2.3.1 Ergebnisse

```
Vorname: klaus
Nachname: maier
Wohnort: meckenbeuren
Geburtsjahr: None
Ort: meckenbeuren
Tätigkeit: industriemechaniker
Hobby: politik
Institution: p+w metallbau gmbh & co. kg
E-Mails: [maier@pw-metallbau.de]
Kontaktinformation: []

Phishing-Mail:
Betreff: Industriemechaniker bei der P+W Metallbau GmbH & Co. KG
Hallo Herr Maier,
als Industriemechaniker bei der Institution P+W Metallbau GmbH & Co. KG,
stehen Ihnen nun alle Möglichkeiten offen. Sehen Sie sich die neuen Mög-
lichkeiten unter folgendem Link an.
https://badlink.com

Mit freundlichen Grüßen

Ihr Karriere-Team der Institution P+W Metallbau GmbH & Co. KG
```

Bild 10.3: Programmausgabe zum Testfall 3 - "Klaus Maier" und "Meckenbeuren"

```
Vorname: klaus
Nachname: maier
Wohnort: tettnang
Geburtsjahr: None
Ort: meckenbeuren
Tätigkeit: bäcker
Hobby: reisen
Institution: europäische fachhochschule
E-Mails: []
Kontaktinformation: []

Phishing-Mail:
Betreff: Bäcker bei der Institution Europäische Fachhochschule
Hallo Herr Maier,
als Bäcker bei der Institution Europäische Fachhochschule, stehen Ihnen
nun alle Möglichkeiten offen. Sehen Sie sich die neuen Möglichkeiten un-
ter folgendem Link an.
https://badlink.com

Mit freundlichen Grüßen

Ihr Karriere-Team der Institution Europäische Fachhochschule
```

Bild 10.4: Programmausgabe zum Testfall 3 - “Klaus Maier“ und “Tettnang“

10.3 Übersicht und Bewertung der erzielten Ergebnisse

10.3.1 Bewertung von Testfall 1

Im Testfall 10.2.1 wurde mit Hilfe eines vollständigen Namens, dem Wohnort und dem einmaligen Instagram-Benutzername gesucht. Dabei wurde ein nahezu richtiges Personenprofil erstellt.

Das Geburtsjahr, der Ort, das Hobby und die Kontaktinformation ist gefunden worden und spricht mit dem tatsächlichen Personenprofil überein. Das Hobby “Fitness“ wird auf dem Instagram-Profil der Ziel- sowie der Kontaktperson gefunden. Allerdings hat der Kontakt keinen vollständigen Namen angegeben. Dadurch konnte nur der Vorname “Sophie“ gefunden werden. Die Tätigkeit ist allerdings nicht korrekt. Hierbei wurde auf

einer Webseite ein Nachname falsch interpretiert und somit als Tätigkeit festgelegt. Das hat den Grund, dass keine weitere Tätigkeit gefunden wurde und somit “Maurer” höchste Gewichtung hat.

Die Anrede für die Phishing-Mail wurde korrekt ausgewählt. Die gewonnene Kontaktinformation ist verwendet worden, um das Opfer zu täuschen. Des Weiteren ergibt die E-Mail Sinn und weist eine korrekte Grammatik auf.

10.3.2 Bewertung von Testfall 2

Der Testfall 10.2.2 zeigt ebenfalls ein annähernd perfektes Ergebnis. Alle Personenattribute bis auf das Hobby sind korrekt und stimmen mit der gesuchten Person überein. Die gefundenen Daten sind Ort, Tätigkeit, Hobby und Institution. Dabei wird die gefundene Tätigkeit und Institution für die Generierung der Phishing-Mail verwendet. Das Hobby “Motorrad” ist nicht richtig. Allerdings wird dieses Hobby auf zwei unterschiedlichen Webseiten gefunden. Dadurch entsteht eine Wertung von 0.34. Das korrekt Hobby “Basketball” wurde ebenfalls gefunden, aber mit dem Score 0.33. Somit wurde laut der Berechnung das richtige Hobby ausgewählt, dennoch ist es in diesem Fall falsch.

10.3.3 Bewertung von Testfall 3

Bei dem Testfall 10.2.3 wurden zwei Personensuchen für das selbe Opfer mit unterschiedlichen Orten durchgeführt. Dabei ist festzustellen, dass das gefundene Profil mit dem Ort “Meckenbeuren“, bis auf den angegebenen Wohnort, vollständig mit der gesuchten Person übereinstimmt. Wogegen der tatsächliche Wohnort zu einem komplett fehlerhaften Profil führt. Grund hierfür ist, dass diese Person in Meckenbeuren arbeitet und alle gefundenen Einträge des Opfers berufsbezogen sind.

Bei diesem Ergebnis ist zu erkennen, wie ausschlaggebend der angegebene Wohnort für die Suche ist.

11 Fazit und Ausblick

11.1 Fazit

Das Ergebnis der Testfälle ist überwiegend positiv. Dennoch ist nur eins der vier Testergebnisse vollständig korrekt. Demnach ist die Antwort auf die Forschungsfrage, ob es möglich ist, ausschließlich korrekte Opferprofile zu erstellen, nein. Das hat den Grund, dass eine große Anzahl an Elementen in den Wortsammlungen dazu führt, dass viele inkorrekte Schlüsselwörter gefunden werden. Somit steigt die Gefahr von Fehlinformationen. Möglicherweise könnte diese Gefahr durch Methoden, welche unter strengeren Bedingungen die Information von Webseiten herausfiltern, gesenkt werden. Dennoch ist die Mehrzahl der Personenattribute in den meisten Fällen richtig. Der Aufwand zu Erstellung einer Phishing-E-Mail ist durch die Automatisierung verschwindend gering. Lediglich die variierende Laufzeit der Anwendung muss beachtet werden. Diese ist abhängig von der gefundenen Information. Somit können glaubwürdige Phishing-Mails mit allen Kriterien erstellt werden. Dennoch hängt die Glaubwürdigkeit und der Erfolg solch einer Mail von dem Charakter eines Opfers ab.

Die definierten Ziele in Kapitel 1.2 sind überwiegend erfüllt. Die erstellte Suchfunktion bietet die Möglichkeit bekannte Daten über die Zielperson einzugeben. Diese Daten dienen zur Identifizierung der gesuchten Person und ermöglichen das Auslesen von bedeutender Information. Allerdings ist eine vollständig korrekte Identifizierung der Zielperson nicht möglich.

E-Mail-Adressen werden aus den Webseiten herausgelesen. Falls keine übereinstimmende Adresse gefunden wird, generiert ein Algorithmus einen Pool an möglichen Adressen. Um zu Beweisen, dass die Phishing-Mail mit der entwickelten Anwendung versendet werden kann, wurde eine Zieladresse festgelegt. Des Weiteren wird der Inhalt einer Mail

abhängig von den gewonnen Informationen ausgewählt und mit den entsprechenden Daten ergänzt.

11.2 Ausblick

Um die Personenidentifikation zu erweitern, könnten Bilderkennungen verwendet werden. Dadurch dienen gleiche Profilbilder auf unterschiedlichen Social-Media-Plattformen als weitere Identifikationsschlüssel. Für diese Methode eignet sich eine Bildererkennungssoftware oder die Google-Bildersuche. Eine weitere Optimierung der Personenidentifizierung kann das Beachten von Zeiträumen sein. Dabei wird erkannt, ob der Inhalt oder das Erstellungsdatum einer Webseite mit dem Alter der Person grundsätzlich übereinstimmt. Hierfür können Jahreszahlen und mögliche Metadaten der Webseite beziehungsweise der Domain ausgelesen werden.

Es stellt sich die Frage, ob die Laufzeit der Anwendung durch die Optimierung der Methode zur Erkennung von wichtigen Informationen verbessert werden kann. Dafür wäre es denkbar, den Vergleich der Schlüsselwörter mit den Elementen der Wortsammlungen zu optimieren. Dazu werden die Datenbanken sortiert und die Schlüsselwörter mit einem angewendeten Suchalgorithmus verglichen. Des Weiteren könnte ein neuronales Netz trainiert werden. Als Trainingsdaten können die Wortsammlungen mit den entsprechenden Kategorien dienen. Das dabei entstehende Netz würde beispielsweise eigenständig das Schlüsselwort "Fußball" aus dem Text herauslesen und in die Kategorie Hobby einordnen. Des Weiteren können die gefundene Elemente durch Stemming auf den Wortstamm zurückgeführt werden. Wodurch eine verbesserte Wertung der vorkommenden Elemente entstehen könnte.

Damit die Wahrscheinlichkeit erhöht wird, dass sich die korrekte E-Mail-Adresse in dem erzeugten Adresspool befindet, können weitere Adresse generiert werden. Als Ideengeber könnte hierfür das OSINT-Tool [Baz19a] dienen. Darüber hinaus können dem Adresspool mögliche Firmenadressen hinzugefügt werden. Dazu müsste allerdings die Institution der Zielperson bekannt sein. Der erzeugt Adresspool beinhaltet viele mögliche E-Mail-Adressen der gesuchten Person. Jedoch ist nicht jede dieser Adressen gültig. Aus diesem Grund

können die generierten Adressen validiert werden. Möglichkeiten dafür sind bereitgestellte Webseiten oder ein eigenes Skript.

Aktuell wird die Phishing-E-Mail mit der Adresse des gefälschten GMX-Accounts versendet. Dadurch steht diese Adresse als Absender in der entsprechenden Mail. Um die Glaubwürdigkeit der Phishing-Mail zu steigern, kann die Absenderadresse verschleiert werden. Dies ist möglich, indem der E-Mail-Header verändert wird. Im Fall, dass Informationen über Kontakte der Zielperson gefunden werden, können diese Daten zur Generierung einer gefälschten Absenderadresse verwendet werden. Darüber hinaus könnte bei einer bekannter Institution die dazugehörige Domain der realen E-Mail-Adresse herausgefunden werden.

A Ausschnitt der Wortsammlung

“Tätigkeit“

Altenpflegehelfer	Baustoffprüfer	Bankkauffrau
Altenpfleger	Bautenbeschichter	Baugeräteführerin
Änderungsschneider	Bauwerksabdichter	Baustoffprüferin
Fachassistentin	Bauwerksmechaniker	Bautenbeschichterin
Anlagenmechaniker	Bauzeichner	Bauwerksabdichterin
Artist	Beamter	Bauwerksmechanikerin
Asphaltbauer	Behälterbauer	Bauzeichnerin
Assistent	Bekleidungsnäher	Beamtin
Stimmlehrer	Bekleidungstechnische	Behälterbauerin
Aufbereitungsmechaniker	Bergmann	Bekleidungsnäherin
Augenoptiker	Bergbautechnologe	Bergbautechnologin
Ausbaufacharbeiter	Berufsflugzeugführer	Berufsflugzeugführerin
Automatenfachmann	Berufshubschrauberführer	Berufshubschrauberführerin
Automobilkaufmann	Berufskraftfahrer	Berufskraftfahrerin
Atemlehrer	Bestattungsfachkraft	Bestattungsfachkraft
Sprechlehrer	Betonbauer	Betonbauerin
Altenpflegehelferin	Betonfertigteilbauer	Betonfertigteilbauerin
Altenpflegerin	Binnenschiffer	Binnenschifferin
Änderungsschneider	Biologielaborant	Biologielaborantin
Fachassistentin	Biologiemodellmacher	Biologiemodellmacherin
Anlagenmechanikerin	Biologisch-technische	Bodenlegerin
Artistin	Bodenleger	Bogenmacherin
Asphaltbauerin	Bogenmacher	Bootsbauerin
Assistentin	Bootsbauer	Böttcherin
Stimmlehrerin	Böttcher	Brauerin
Aufbereitungsmechanikerin	Brauer	Brennerin
Augenoptikerin	Brenner	Brunnenbauerin
Ausbaufacharbeiterin	Brunnenbauer	Buchbinderin
Automatenfachfrau	Buchbinder	Buchhändlerin
Automobilkauffrau	Buchhändler	Büchsenmacherin
Atemlehrerin	Büchsenmacher	Bühnenmalerin
Sprechlehrerin	Bühnenmaler	Bühnentänzerin
Bäcker	Bühnentänzer	Bürstenmacherin
Bankkaufmann	Bürstenmacher	Chemielaborant
Baugeräteführer	Bäckerin	Chemielaborjungwerker

Bildverzeichnis

2.1	Architektur eines Web Crawlers [Cas08]	9
8.1	Aufbau des entwickelnden Web Crawlers	33
8.2	Container eines Google-Suchergebnisses [LLC19]	37
8.3	Beispieltext [SS11]	45
8.4	Schlüsselwörter mit zugehörigem Score	47
8.5	Funktion des Algorithmus zur Suche nach dem Geburtsjahr	50
8.6	Container mit Profil-Vorschlägen [INS19]	57
8.7	Pop-up-Fenster mit abonnierten Profilen [INS19]	58
8.8	Personenklasse	61
9.1	Ein berufliches E-Mail-Muster mit den Personenattributen Nachname und Institution	66
9.2	Ein berufliches E-Mail-Muster mit den Personenattributen Nachname und Tätigkeit	67
9.3	Ein berufliches E-Mail-Muster mit dem Personenattribut Nachname und einer festgesetzten Tätigkeit	67
9.4	Ein berufliches E-Mail-Muster mit den Personenattributen Nachname, Tätigkeit und Institution	67
9.5	Ein privates E-Mail-Muster mit den Personenattributen Vorname, Kontaktname und Kontaktinformation	68
9.6	Ein privates E-Mail-Muster mit den Personenattributen Vorname und Hobby	69
9.7	Ein privates E-Mail-Muster mit den Personenattributen Vorname und Jahrgang	69
9.8	Ein privates E-Mail-Muster mit den Personenattributen Vorname und Ort	69
10.1	Programmausgabe zum Testfall 1 - "Marco Lang", "Tettnang" und "Instagram-Benutzername"	74
10.2	Programmausgabe zum Testfall 2 - "Anita Schmidt" und "Heidesheim"	75
10.3	Programmausgabe zum Testfall 3 - "Klaus Maier" und "Meckenbeuren"	76
10.4	Programmausgabe zum Testfall 3 - "Klaus Maier" und "Tettnang"	77

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.
- [Baz18] BAZZELL, MICHAEL: *Open Source Intelligence Techniques: Resources for Searching and Analyzing Online Information*. CreateSpace Independent Publishing Platform, USA, 6th , 2018.
- [Baz19a] BAZZELL, MICHAEL: *Email Assumptions*. <https://inteltechniques.com/osint/email.html>, 2019. Abrufdatum: 01.02.2019.
- [Baz19b] BAZZELL, MICHAEL: *Intel Techniques*. <https://inteltechniques.com/index.html>, 2019. Abrufdatum: 20.04.2019.
- [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.
- [Bun] BUNDESAMT FR SICHERHEIT IN DER INFORMATIONSTECHNIK: *Aktuelle Information zur Schadsoftware Emotet*. <https://www.bsi-fuer-buerger.de/BSIFB/DE/Service/Aktuell/Informationen/Artikel/emotet.html?nn=6775642>. Abrufdatum: 22.04.2019.
- [Cas08] CASTILLO, CARLOS: *File:WebCrawlerArchitecture.svg*. <https://commons.wikimedia.org/wiki/File:WebCrawlerArchitecture.svg>, 2008. Abrufdatum: 25.04.2019.
- [CH15] CHRISTOPHER HADNAGY, MICHELE FINCHER: *Phishing Dark Waters: The Offensive and Defensive Sides of Malicious E-mails*. 2015.
- [CT⁺94] CAVNAR, WILLIAM B, JOHN M TRENKLE.: *N-gram-based text categorization. Proceedings of SDAIR-94, 3rd annual symposium on document analysis and information retrieval*, 161175. Citeseer, 1994.
- [DAT15] DATEV UND DEUTSCHLAND SICHER IM NETZ E.V.: *Verhaltensregeln zum Thema "Social Engineering"*. 2015.
- [Dat19] DATAREPORTAL, WE ARE SOCIAL; HOOTSUITE,: *Ranking der gren sozialen Netzwerke und Messenger nach der Anzahl der monatlich aktiven Nutzer (MAU) im Januar 2019 (in Millionen)*. <https://de.statista.com/statistik/daten/studie/181086/umfrage/die-weltweit-groessten-social-networks-nach-anzahl-der-user/>, 2019. Abrufdatum: 24.04.2019.
- [dev18] DEVELOPERS, SCRAPY: *Scrapy at a glance*. <http://doc.scrapy.org/en/latest/intro/overview.html>, 2018. Abrufdatum: 28.02.2019.
- [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.
- [Gmb18] GMBH, BIBLIOGRAPHISCHES INSTITUT: *Daten*. <https://www.duden.de/suchen/dudenonline/daten/>, 2018. Abrufdatum: 22.04.2019.
- [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.
- [INS19] INSTAGRAM, INC.: *Instagram*. <https://www.instagram.com/>, 2019. Abrufdatum: 25.04.2019.
- [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.
- [Pro15] PROEBSTER, WALTER: *Rechnernetze: Technik, Protokolle, Systeme, Anwendungen*. Walter de Gruyter GmbH & Co KG, 2015.

- [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.
- [Sei69] SEIBICKE, WILFRIED: *Wie schreibt man gutes Deutsch?: eine Stilfibel*. Bibliogr. Inst., 1969.
- [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.kidscodecs.com/stop-words/>. Abrufdatum: 29.01.2019.
- [SS11] SCHUBERT, SIGRID ANDREAS SCHWILL: *Didaktik der Informatik. Didaktik der Informatik*, 1–30. Springer, 2011.
- [The01] THELWALL, MIKE: *A web crawler design for data mining*. Journal of Information Science, 27(5):319–325, 2001.
- [W3S] W3SCHOOLS: *HTML URL Encoding Reference*. https://www.w3schools.com/tags/ref_urlencode.asp. Abrufdatum: 27.02.2019.