

Erstellung eines Programms zur automatisierten Informationsbeschaffung von personenbezogenen Daten in Verbindung mit einem automatisierten Phishing-Mailgenerators

Bachelorarbeit

Social Engineering

im Studiengang **Angewandte Informatik**

an der Hochschule Ravensburg - Weingarten

von

Marco Lang **Matr.-Nr.: 27416**

Abgabedatum : 10. Februar 2019

Erklärung

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

Generierung eines personalisierten Mail-Generators

selbstständig angefertigt, nicht anderweitig zu Prüfungs Zwecken vorgelegt, keine anderen als die angegebenen Hilfsmittel benutzt und wortliche sowie sinn-gemaesse Zitate als solche gekennzeichnet habe.

Weingarten, 10. Februar 2019

Autor Name

Inhaltsverzeichnis

Kurzfassung	IV
Abstract	V
Danksagung	VI
Vorwort	VII
1 Einleitung	1
1.1 Motivation	1
1.2 Zielsetzung	1
1.3 Eigene Leistung	3
1.4 Aufbau der Arbeit	3
2 Grundlagen	4
2.1 Social Engineering	4
2.1.1 Phishing	5
2.2 Informationsbeschaffung im Internet	7
2.2.1 Web Scraping	7
2.2.2 Web Crawler	7
2.3 Personenbezogene Daten	7
2.4 Natural Language Processing	8
2.5 Textanalyse	8
2.5.1 Stoppwörter	8
3 Problemspezifikation	9
4 Anforderungsanalyse und Priorisierung	10
4.1 Anforderung an das Programm bzw. an die Programmiersprache	10
4.2 Anforderung an die Informationsbeschaffung	10
4.2.1 Informationsbeschaffung von einer ausgewählten Person	11
4.2.2 Informationsbeschaffung von unbestimmten Personen	11
4.3 Anforderung an die Datenverwaltung/-speicherung	11

4.4	Anforderung an die Generierung der E-Mail-Adressen	12
4.5	Anforderung an die E-Mail-Muster	12
4.6	Anforderung an die Erstellung der Phishing-Mail	12
4.7	Weitere Anforderungen	13
4.8	Priorisierung	13
5	Lösungsideen	14
5.1	Programmiersprache/ GUI	14
5.2	Informationsbeschaffung einer ausgewählten Person	14
5.2.1	Wie sieht die Suche nach einer Person im Internet aus?	14
5.2.2	Wann handelt es sich um die gesuchte Person?	16
5.2.3	Wie wird wichtige Information auf einer Website erkannt?	18
5.2.4	Speicherung der gewonnenen Daten	22
5.3	Informationsbeschaffung von einer großen Menge unbekannter Personen .	22
5.4	Generierung der E-Mail-Adressen	23
5.5	Erstellung der E-Mail-Muster	24
5.6	Erzeugung der Phishing-Mail	24
6	Bewertung der Lösungsideen anhand der Anforderung	25
7	Umsetzung	27
7.1	Textanalyse mit Hilfe von Python NTLK	27
7.2	Informationsbeschaffung von der Website www.fupa.net	28
7.2.1	Erstellung eines Web Crawlers	28
7.3	Datenverwaltung und Speicherung	29
7.3.1	Speicherung von Personendaten in CSV oder mySQL	29
8	Evaluation der Implementation	30
10	Hauptteil	32
10.1	Hauptteil	32
10	Hauptteil	32
10.1	Hauptteil	32
11	Ethische und rechtliche Betrachtung	33
11.1	Ethische Betrachtung	33
12	Schlussbemerkungen und Ausblick	34
A	Ein Kapitel des Anhangs	35

Glossar	36
Abkürzungsverzeichnis	37
Symbolverzeichnis	38
Literatur	39
Stichwortverzeichnis	41

Kurzfassung

Abstract

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

Danksagung

Vorwort

1 Einleitung

1.1 Motivation

Laut dem Bundeskriminalamt hat sich die Zahl der Cyberkriminalität mit einem klaren Trend nach oben entwickelt. [Bun18] Aus diesem Grund werden System immer sicherer und Firewalls immer noch besser. Das hat zu Folge, dass Angreifer oft auf Methoden ausweichen, bei denen der Mensch als Schwachstelle des Systems ausgenutzt wird. Daher ist eine häufig verwendete Technik von Cyberkriminalität das E-Mail-Phishing.

In den neusten Fällen von Phishing-Mail-Attacken zeigt die Verbraucherzentrale Nordrhein-Westfalen, dass diese meist direkt an eine Person adressiert sind. Das heißt, in dieser Art von E-Mail, werden personenbezogene Daten verwendet. Ein Beispiel dafür, sind die gefälschten DSGVO-E-Mails. Hier wird die Zielperson im Namen der Sparkasse, persönlich mit Namen angesprochen. [NW18]

Solch ein Angriff benötigt im Voraus eine ausführliche Recherche über das Opfer. Als Informationsquelle für die Recherche können beliebig viele Quellen verwendet werden. Jedoch ist in der heutigen Zeit das Internet eine der meistgenutzten Informationsquellen. [All18]

1.2 Zielsetzung

Ziel dieser Arbeit ist es ein Programm zu entwickeln, welches automatisiert nach personenbezogenen Daten im Internet sucht und daraus eine Phishing-Mail generiert. Dabei soll der Fokus auf der automatisierten Informationsbeschaffung liegen.

Es sollen grundsätzlich zwei verschiedene Suchfunktionen mit diesem Programm möglich

sein.

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

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

Ziel 2 *Webseiten, die eine große Menge von personenbezogener Daten enthalten, auslesen und analysieren.*

Durch die zweite Suchfunktion soll eine große Menge an Daten gewonnen werden und dadurch ein weitläufiger Angriff zu simulieren.

Bei der zweiten Suchfunktion sollen nur bestimmte Webseiten vorgegeben werden, welche ausgelesen und analysiert werden sollen. Durch diese Funktion ist es möglich einen weitläufigen “*real-world*“ Phishing-Mail-Angriff zu simulieren.

Ziel 3 *E-Mail-Adressen aus den gewonnenen Daten generieren.*

Durch die Zusammensetzung von Vorname, Name und Geburtsjahr und/oder Firma werden die E-Mail-Adressen generiert.

Ziel 4 *Phishing-Mail-Muster erstellt*

Abhängig von den gefundenen Informationen, soll mit Hilfe der Muster eine Phishing-Mail mit glaubhaftem und sinnvollem Inhalt erstellt werden.

Ziel 5 *Phishing-Mail erzeugen.*

Mit der vorhandenen Information, der E-Mail-Adresse und einem passende Muster, soll eine Phishing-Mail erzeugt und versendet werden können.

1.3 Eigene Leistung

In dieser Arbeit wird ein Programm erstellt, welches personenbezogene Daten automatisiert aus dem Internet heraussucht und diese in potentielle Opferprofile ablegt. Die gewonnenen Informationen werden automatisiert in eine personalisierte Phishing-E-Mail eingebaut. Für einen höheren Erfolg werden E-Mail-Muster erstellt.

Damit ein kompletter Ablauf eines Phishing-Mail-Angriffs simuliert werden kann, wird ein Algorithmus entwickelt, der aus den gewonnen Informationen eine E-Mail-Adresse generiert.

1.4 Aufbau der Arbeit

Die Arbeit gliedert sich in einem theoretischen und praktischen Teil auf. Der Theorie-Teil beginnt im zweiten Kapitel und beschreibt die Grundlagen im Bereich von Social Engineering, der Informationsbeschaffung im Internet, personenbezogene Daten und der Textanalyse. Im nächsten Kapitel befindet sich die Anforderungsanalyse, bei der die Anforderungen an die Arbeit festgelegt werden. Darauf folgen die Lösungsvorschläge im Kapitel vier und die Auswahl der Lösung anhand den Anforderungen im Kapitel fünf. Anschließend wird bei der Umsetzung auf den Praktischen Teil eingegangen. Am Ende befindet sich das Fazit, der Ausblick und der Anhang.

2 Grundlagen

2.1 Social Engineering

Die Definition von Social Engineering (SE) ist nicht eindeutig. Es gibt sehr verschiedene Ansichten von der Definition. Die Idee von Social Engineering ist, eine Ziel so zu manipulieren, damit das Ziel eine für den Angreifer bessere Entscheidung trifft. In dem Buch Social Engineering - The Art of Human Hacking, von Christopher Hadnagy, ist Social Engineering definiert als:

“social engineering is the act of manipulating a person to take an action that may or may not be in the “target’s” best interest“ [Had11]

Wiederum lautet die Definition in dem Buch von Kevin D. Mitnick:

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

SE 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, o.ä. bekommt. Darüber hinaus ist SE in vielen Berufen ein täglicher Bestandteil. Beispielsweise manipulieren Ärzte viele Patienten mit einer Placebo-Behandlung. Bei dieser Behandlung wird dem Patient ein wirkstoff-freies Medikament verschrieben. Ausschließlich durch die Manipulation des Patienten und den sogenannten Palzebo-Effekt können Erfolge erzielt werden.

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

Der Aufbau eines SE-Angriffes ist definiert in mehrere Phasen. Das wohl bekannteste Modell für einen Social Engineering-Angriffszyklus ist in dem Buch von Kevin D. Mitnick's [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 viel 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 *Exploiting Trust-Phase* ausgenutzt. Hier will der Angreifer die eigentlich Information vom Opfer herausfinden. Dies geschieht einerseits durch bestimmtes Nachfragen oder Manipulation. *Utilize Information* ist die letzte Phase. Dort wird die gewonnene Information genutzt um das eigentliche Ziel des Angreifers zu erreichen.

Grundsätzlich werden bei einem Social Engineering Angriff menschliche Wünsche, Ängste und verbreitete Verhaltensmuster verwendet um ein Opfer zu manipulieren. [uDsiNe15]

2.1.1 Phishing

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

Die wohl bekannteste Angriffsmethode von Phishing ist das E-Mail-Phishing. Bei diesem Verfahren, versendet ein Angreifer meist eine gefälschte E-Mail, um ein Opfer zu täuschen und dadurch sein Ziel zu erreichen. Die sogenannten Phishing-Mails enthalten meist eine Aufforderung einen Link zu öffnen und sehen täuschend echt aus. Zum Beispiel könnte der Angreifer ein Layout von Amazon verwenden und das Ziel auffordern, den Link zu

öffnen wegen einem Authentifizierungsproblem. Nachdem Sie auf den Link geklickt haben müssen Sie sich anmelden. Hier könnten die Angreifer Ihre Anmeldedaten abgreifen, nachdem sie das Opfer eingeben hat. Sobald die Anmeldedaten eingegeben wurden, könnte eine Fehlermeldung erscheinen, die sagt: "Hoppla, ein Fehler ist aufgetreten, melden Sie sich bitte neu an!". Anschließend wird die originale Seite geladen, das Opfer kann sich korrekt anmelden und der Angreifer hat ohne einen großen Aufwand die Anmeldedaten der Zielperson.

Für diese Methode benötigt der Angreifer nicht nur Social Engineering sondern auch technische Fähigkeiten. [CH15]

Spear-Phishing

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

Bei dieser Form von E-Mail-Phishing spielt die Opferauswahl und die Informationsbeschaffung eine wichtig 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 Informationsbeschaffung wird oft in Kauf genommen, da der Erfolg bei diese Methode vielversprechender ist al beim herkömmlichen E-Mail-Phishing. 91% der Advanced Persistent Threat (APT) Angriffe auf Firmen beginnen mit einer Spear-Phishing-E-Mail. Die Schadsoftware wir meisten als Remote Access Trojans (RATs) in einem Zip-Datei überliefert. [Cal13]

2.2 Informationsbeschaffung im Internet

2.2.1 Web Scraping

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

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

2.2.2 Web Crawler

Web Crawler sind Computerprogramme, die mit Hilfe der Hypertextstruktur das Internet durchlaufen. 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]

2.3 Personenbezogene Daten

Laut dem DSGVO sind *personenbezogene Daten* “alle Informationen, die sich auf eine identifizierte oder identifizierbare natürliche Person (im Folgenden „betroffene Person“) beziehen; als identifizierbar wird eine natürliche Person angesehen, die direkt oder indirekt, insbesondere mittels Zuordnung zu einer Kennung wie einem Namen, zu einer Kennnummer, zu Standortdaten, zu einer Online-Kennung oder zu einem oder mehreren

besonderen Merkmalen identifiziert werden kann, die Ausdruck der physischen, physiologischen, genetischen, psychischen, wirtschaftlichen, kulturellen oder sozialen Identität dieser natürlichen Person sind;“ [DSG]

2.4 Natural Language Processing

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

Stemming ist eine Methode der Wortstandardisierung, bei der verwandte Wörter auf ihrer Stammform reduziert werden. Dabei wird bei dem Rechengang auf den Stamm und die Semantik eines Wortes geachtet. Aus diesem Grund fällt der Name Stammformreduktion öfters in Verbindung von Stemming. [EAD09]

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

2.5 Textanalyse

2.5.1 Stoppwörter

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]

3 Problemspezifikation

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

Dass hier beschriebene Problem zeigt, dass der Zugang für persönliche Information durch das Internet für die Öffentlichkeit einfacher gemacht wird. Es soll mit einem kritisch Blick darauf gezeigt werden, mit welchem Aufwand, personenbezogene Daten aus dem Internet herausgelesen, analysiert und für einen Phishing-Mail-Angriff verwendet werden kann.

4 Anforderungsanalyse und Priorisierung

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

4.1 Anforderung an das Programm bzw. an die Programmiersprache

Es soll eine möglichst übersichtliche und performante Skriptsprache verwendet werden, mit der eine automatisierte Informationsbeschaffung gut möglich ist. Eine Eingabe über die Konsole oder über eine graphische Benutzeroberfläche soll ebenfalls möglich sein. Aus diesem muss die Programmiersprache keine GUI-Programmierung mit sich bringen.

4.2 Anforderung an die Informationsbeschaffung

Die Anforderung an die Informationsbeschaffung von personenbezogenen Daten lässt sich in zwei Teile gliedern. Der erste Teil beinhaltet die Informationsbeschaffung von ausgewählten Personen und der zweite Teil die Informationsbeschaffung von einer großen Menge unbekannten Personen.

4.2.1 Informationsbeschaffung von 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. Dies soll mit Hilfe eines *Web-Crawlers* und mit einem *Web-Scraper* umgesetzt werden.

Das zu entwickelnde Programm soll für die Suche bekannte Daten wie Vorname, Nachname, Geburtsjahr, Ort und Benutzernamen von Social Media Plattformen über eine Konsolen oder eine grafische Oberfläche einlesen können.

Die Herausforderung besteht darin, zu erkennen, wann und ob es sich um die Information der gesuchten Person handelt.

4.2.2 Informationsbeschaffung von unbestimmten Personen

Es soll eine Prototyp-Suchfunktion entwickelt werden, die eine komplette Website nach personenbezogenen Daten durchsucht. Dabei sollen möglichst viele Informationen von möglichst vielen Personen herausgefunden werden. Jedoch sind diese Personen dem Programm-Anwender unbekannt. Die Informationen werden aus Webseiten mit einer großen Anzahl von Mitgliedern herausgelesen. Bei dieser Suchfunktion soll den Anwender aus den vorgegebenen Webseiten eine Seite auswählen können. Die ausgewählte Webseite wird daraufhin komplett ausgelesen und nach personenbezogenen Daten durchsucht.

Dabei soll der zu entwickelnde *Web Scraper* möglichst performant arbeiten und kann *hartkodiert* werden. Allerdings müssen E-Mail-Adressen ebenfalls gefunden werden können, obwohl die Position einer E-Mail-Adressen auf einer Webseite variieren kann.

4.3 Anforderung an die Datenverwaltung/-speicherung

Ausgelesene Daten sollen vor dem speichern formatiert und klassifiziert werden, damit die Daten später korrekt in die Phishing-Mails eingesetzt werden können. Die Schwierigkeit besteht darin, zu erkennen, um welche Art von Information es sich handelt. Zusätzlich

sollen die Daten in einer gut übersichtlichen Struktur gespeichert werden und müssen beliebig erweiterbar sein.

4.4 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 soll eine größere Anzahl von möglichen E-Mail-Adressen erzeugt werden. Durch den Pool an erzeugten E-Mail-Adressen soll die Wahrscheinlichkeit erhöht werden, dass die richtige E-Mail-Adresse dabei ist. Des Weiteren sollen die Adresse auf Verfügbarkeit und Gültigkeit geprüft werden.

4.5 Anforderung an die E-Mail-Muster

Bei der Erstellung der E-Mail-Muster handelt es sich ausschließlich um das Erstellen potentieller Inhalte einer E-Mail, welcher mit den gewonnenen Information über eine Person erweitert werden kann. Die Muster sollen erstellt werden und so klassifiziert sein, dass für jedes gefundene Opferprofil ein passendes Muster vorhanden ist. Des Weiteren soll der E-Mail-Text mit den eingesetzten Informationen Sinn ergeben und eine korrekte Grammatik beinhalten. Weiterführend können SE-Fähigkeiten genutzt werden um die Zielperson tatsächlich zu manipulieren und täuschen. Hierfür können beispielsweise Gefühle wie Freude und Angst ausgenützt oder gefälschte E-Mails von bekannten Firmen in Betracht gezogen werden.

4.6 Anforderung an die Erstellung der Phishing-Mail

Die Phishing-Mails sollen automatisiert erstellt werden. Die Auswahl des richtigen E-Mail-Musters zu der gewonnenen Opferinformation soll ebenfalls automatisiert ablaufen.

4.7 Weitere Anforderungen

Unter anderem soll die Arbeit Antworten auf die folgenden Fragen finden. Mit welchem Aufwand ist eine Phishing-Mail-Angriff verbunden? Ist es möglich ein Personenprofil zu erstellen, bei dem ausschließlich korrekte Informationen vorhanden sind?

4.8 Priorisierung

Die Tabelle 4.1 zeigt die Priorisierung der Anforderungen. Dabei liegt der eindeutige Fokus auf der Informationsbeschaffung von personenbezogenen Daten und der Erstellung von E-Mail-Mustern.

Tabelle 4.1: Priorisierung der Anforderungen

Anforderung	Priorisierung (A-C)
Informationsbeschaffung von ausgewählten Personen	<i>A</i>
Informationsbeschaffung von vielen unbekannten Personen	<i>A</i>
E-Mail-Muster erstellen	<i>A</i>
Phishing-Mail erzeugen	<i>B</i>
Datenverwaltung/-speicherung	<i>B</i>

5 Lösungsideen

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

5.1 Programmiersprache/ GUI

Für die Auswahl der Programmiersprache gibt es viele Auswahlmöglichkeiten. Dennoch wird in dieser Abschlussarbeit die Programmiersprache Python verwendet, da sie die Nötigen Eigenschaften mit sich bringt.

Für die Eingabe von Suchdaten, besteht für beide Informationsbeschaffungen die Möglichkeit eine Grafische-Bedienoberfläche oder Konsolen-Eingabe zu verwenden.

5.2 Informationsbeschaffung einer ausgewählten Person

5.2.1 Wie sieht die Suche nach einer Person im Internet aus?

Die Suche nach einer Person im Internet kann durch mehrere Ansätze erfolgen. Die nachstehenden Ansätze unterscheiden sich in der Art Suche und in dem Umgang der eingegeben Daten.

Die Art der Personensuche wird anhand den eingegebenen Daten angepasst

Abhängig von der Anzahl und Art der Daten, die von dem Programm-Anwender eingegeben wurden, wird die Art und Reihenfolge der Suche variiert. Die nachfolgenden Fälle sollen diesen Ansatz verdeutlichen.

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

Ein weiterer Fall beschreibt das Szenario, wenn ein Benutzername von der gesuchten Person in das Programm eingegeben wird. Hierbei handelt es sich um einen Benutzernamen von Webseiten wie Facebook, Instagram, usw. Zuallererst, kann hier die entsprechende Webseite nach Informationen zu dem angegebenen Benutzername durchsucht werden. Dadurch können zusätzliche Daten herausgefunden werden, die bei der weiteren Suche von Vorteil wären. Nachdem die Webseite nach dem Nutzernamen durchsucht und ausgewertet wurde, kann nun mit herkömmlichen Suchmaschinen die Suche erweitert werden.

Es wird unabhängig von den eingegebenen Daten direkt mit einer Suchmaschine nach der Person gesucht

Bei diesem Lösungsansatz werden ausschließlich die herkömmlichen Suchmaschinen verwendet. Die Funktion der Suche besteht darin, dass das Programm den vorgeschlagenen Links der Suchmaschinen folgt, wobei die eingegebenen Daten die Art der Suche nicht beeinflussen.

Nur ausgewählte Webseiten werden nach einer Person durchsucht

Unabhängig von den eingegebenen Daten, werden verschiedene Webseiten durchsucht. Allerdings ohne die Verwendung einer Suchmaschine. Vorschläge für die ausgewählten Webseiten sind Facebook, FuPa, Instagram, Xing, LinkedIn und Twitter.

5.2.2 Wann handelt es sich um die gesuchte Person?

Bei jeder einzelnen Suchvariante, 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. Um diesem Problem entgegen zu wirken, kann die Art der Suche anhand den eingegebene Daten angepasst werden. Dies entspricht dem Ansatz 5.2.1. Die Suche kann dadurch verfeinert werden und die Anzahl der fehlerhaften Vorschläge wird geringer. Dadurch wird die Wahrscheinlichkeit höher, dass es sich um die richtige Person handelt.

Darüber hinaus kann die Personensuche mit einer Suchmaschine durch verbesserte Suchbefehle ebenfalls verfeinert werden. In dem Buch *“Open Source Intelligence Techniques”* [Baz18], werden Suchbefehle für bekannte Suchmaschinen aufgezeigt, mit denen die Suche verbessert werden kann. Dies bedeutet, bei einer Personensuche ist es mit den richtigen Suchbefehlen möglich, die Anzahl der Vorschläge um einen großen Teil zu verringern. Ein Beispiel in dem Buch von Michael Bazzell zeigt, wie es funktioniert von 8770 Vorschlägen auf lediglich neun Vorschläge zu reduzieren. [Baz18] Auch bei dieser Lösungsidee wird die Wahrscheinlichkeit erhöht, dass es sich um die gesuchte Person handelt.

Im Fall dass nach diese Maßnahmen dennoch verschiedene Profile angezeigt werden, können die folgenden Erweiterungen in die Suche mit einfließen.

Erweiterte Kriterien

Hierbei handelt es sich um weitere Kriterien, welche die Suche noch mehr eingrenzen sollen. Bekannte Informationen zur Person sollen dazu dienen, die vorgeschlagenen Seiten

einer Suchmaschine weiter zu filtern. Genau genommen heißt das, dass das Programm in erster Linie nur die Webseite als Informationsquelle verwendet, die alle Suchbegriffe beinhaltet. Darüber hinaus kann das genaue oder grobe Alter der Zielperson mit in die Suche mit einfließen. Dadurch kann erkannt werden ob der Zeitrahmen des Artikels oder das Erstellungsdatum einer Webseite mit dem Alter der Person übereinstimmt.

Kontakte der Suchperson werden in Betracht gezogen

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. In diesem Fall könnten Facebook-Freunden, FuPa-Teammitglieder, Instagram-Follower oder LinkedIn/Xing-Kontakte als Kontaktquelle dienen. Dadurch können weitere Informationen gewonnen werden, die zur Unterscheidung von Profilen nützlich sein könnten.

Profilbilder vergleichen

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

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

Identifikationsschlüssel verwenden

Bekannte Information zur Person können als Identifikationsschlüssel verwendet werden. Allerdings müssen dies einzigartige Daten sein. Dazu zählt die E-Mail-Adresse oder eine Verbindung von mehreren personenbezogene Daten, da der vollständige Name nicht einzigartig ist und sich auf verschiedenen Webseiten unterscheiden kann. Das heißt eine Zielperson kann auf einer Seite einen erfundenen Spitzname und auf der nächsten Seite den

vollständigen Namen verwenden. Bei einer einfachen Suche würde das zu Problemen führen.

Im Fall das auch mit diese Maßnahmen nicht die gesuchte Person identifiziert werden kann, können mehrere Personenprofile erstellt und angezeigt werde. Der Programm-Anwender kann anschließend aus den vorgeschlagenen Profilen eines auswählen.

5.2.3 Wie wird wichtige Information auf einer Website erkannt?

Für die Suche 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 eine *Hartkodierung* zu verwenden, um festgelegte Bereiche einer Webseite auszulesen.

Die Grundidee zur Lösung diese Problems 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 die meisten Sonderzeichen herausgefiltert werden. Sonderzeichen wie "." und "@" werden dabei nicht herausgefiltert, da sie für die E-Mail-Erkennung wichtig sind. Anschließend werden Schlüsselwörter aus dem formatierten Webseitentext generiert. Möglichkeiten zur automatisierten Schlüsselwortgenerierung sind die Verfahren RAKE 5.2.3 und die automatisierte Schlüsselwortgenerierung mit Hilfe von Machine Learning 5.2.3, welche im Laufe dieser Arbeit detailliert beschrieben werden.

Nachdem die Schlüsselwörter generiert und in Listen gespeichert wurden, werden Wortsammlung erstellt. Diese Wortsammlungen sind Listen, welche aussagekräftige Schlüsselwörter enthalten und nach Themen kategorisiert werden. Beispiele für den Inhalt der Listen sind alle Hochschulen und Universitäten in Deutschland, Berufsbezeichnungen/Tätigkeiten, Studiengänge, Hobbys, Städte und Gemeinden in Deutschland.

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

Die Wortsammlungen werden mit Hilfe von bekannter Listen im Internet eigenständig befüllt. Als Informationsquelle dafür dient jegliche Art von Webseite, die nützliche Information enthält. Dazu zählt auch die Seite *Wikipedia - Die freie Enzyklopädie*.

RAKE

RAKE steht für *Rapid Automatic Keyword Extraction* und stellt eine sehr effiziente Methode zur Schlüsselwortgenerierung dar. Die Funktion von RAKE basiert darin, dass Schlüsselwörter mehrere Wörter mit inhaltlicher Relevanz enthalten können, allerdings selten Stoppwörter 2.5.1 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 5.1: Beispieltext

Zu Beginn wird der zu analysierende Text, hier der Beispieltext in Bild 5.1, durch einen Worttrenner in ein Array, bestehend aus möglichen Schlüsselwörtern, aufgeteilt. Das erzeugte Array wird anschließend in Sequenzen von zusammenhängenden Wörtern unterteilt. Dabei erhalten die Wörter in einer Sequenz die gleiche Position und Reihenfolge wie im Ursprungstext und dienen gemeinsam als Kandidatenschlüsselwort. [RECC10]

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

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

Dabei beschreibt der Grad eines Wortes, dass gemeinsame Auftreten mit sich selbst und anderen Schlüsselwörtern. In der Tabelle 5.2.3 ist der Grad für jedes Wort ablesbar, in dem 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“. $deg(w)$

Die Häufigkeit des Vorkommens eines Wortes lässt sich ebenfalls in der Tabelle 5.2.3 ablesen. Allerdings muss hier in der Reihe und Spalte des jeweiligen Wortes nachgeschaut werden. Für das Wort “Wissenschaft” beträgt die Häufigkeit des Vorkommens den Wert 3. Zusammenfassend kann gesagt werden, dass $deg(w)$ die Kandidaten bevorzugt, welche oft und in langen Schlüsselwörtern, die mehrere Wörter enthalten, vorkommen. Dies bedeutet, dass beispielsweise $deg(etabliert)$ eine höhere Bewertung als $deg(informatik)$ bekommt, obwohl beide Wörter gleich oft im Text vorkommen. Dagegen wird bei $freq(w)$, ausschließlich die Häufigkeit des Vorkommens bewertet. Bei der Formel $deg(w)/freq(w)$ werden die Wörter bevorzugt, welche überwiegend in langen Kandidatenwörtern vorkommen. Diese Berechnung 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 höher ist 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 5.2 befindenden Kandidaten mit den dazugehörigen Endbewertungen. [RECC10]

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

Tabelle 5.1: Co-occurrence

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

Bild 5.2: Schlüsselwörter mit zugehörigem Score

Automatic Keyword Extraction mit NLP

Bei dieser Methode wird der vorliegende Text in die einzelnen Wörter unterteilt. Dabei wird eine Liste mit potentiellen Schlüsselwörtern erstellt, in der *Stoppwörter* und Sonderzeichen herausgefiltert werden. Bei den Schlüsselwörtern handelt es sich nicht ausschließlich um ein Wort sondern auch um Wortsequenzen.

Mit Hilfe von Stemming kann nun die Anzahl der Wörter in der Liste weiter reduziert werden, wodurch eine bessere Schlüsselwortgenerierung möglich ist.

Die Liste mit den möglichen Schlüsselwörtern, kann nach der Häufigkeit des Vorkommens eines Wortes im Text sortiert werden. Das bedeutet, dass das Stammwort, welches am Häufigsten im Text vorkommt, in den folgenden Schritten zuerst verwendet wird. Dadurch kann die Laufzeit der Anwendung verbessert werden.

Mit weiteren Regeln, wie eine Mindestanzahl von Buchstaben in einem Wort, können die Schlüsselwörter weiter begrenzt werden.

Keyword Extraction mit Hilfe von Machine Learning

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

5.2.4 Speicherung der gewonnenen Daten

Die gewonnenen Daten können in einem beliebig erweiterbaren Personen-Objekt gespeichert werden. Darüber hinaus lässt sich das Objekt mit bekannten Kontakten der zu suchenden Person erweitern.

Eine andere Möglichkeit wäre die Daten in eine Datei auszulagern. Hierfür wäre eine Datei mit dem Format *CSV* oder *TXT* möglich.

5.3 Informationsbeschaffung von einer großen Menge unbekannter Personen

Für die *real-world* Simulation eines Phishing-Mail-Angriffs werden Webseiten mit großen Menge von personenbezogenen Daten benötigt. Möglichkeiten, ausgenommen von den bekannten Social Media Seiten, sind das Fußballportal FuPa, Xing und LinkedIn.

Zum Auslesen der Webseite kann ein Web Scraper erstellt werden, der es ermöglicht die große Menge von Daten auszulesen. Dieser könnte für die entsprechenden Webseiten

hartkodiert werden. Eine weitere Möglichkeit wäre die Analyse des Webseitentextes, wie bei der Suchfunktion einer ausgewählten Person.

Für die Speicherung der gewonnenen Daten kann eine SQL-Datenbank erstellt werden. Als Alternative kann eine Datei angelegt werden, bei der alle Daten zu allen Personen gut strukturiert gespeichert werden können. Eine Möglichkeit dafür ist das Dateiformat *CSV* oder *TXT*.

5.4 Generierung der E-Mail-Adressen

Eine Möglichkeit zur Generierung der E-Mail-Adressen kann das Open Source-Tool von Michael Bazzell [Baz] sein, welches mit Hilfe eines automatisierten Webbrowsers verwendet werden kann. Bei diesem Tool werden zuerst über ein Formular, Daten für die E-Mail-Generierung eingetragen. Unter anderem sind das Vorname, Nachname und der E-Mail-Provider. Daraufhin werden die vorgeschlagenen E-Mail-Adressen angezeigt, kopiert und in ein Suchfeld eingefügt. Anschließend kann bei Google, Bing, und Facebook nach Einträgen gesucht und falls ein Eintrag gefunden wurde auch angezeigt werden.

Eine Weitere Möglichkeit wäre ein Algorithmus zu entwickeln, der alle möglichen E-Mail-Adressen aus den Kombinationen von Vorname, Nachname, Geburtsjahr, Benutzernamen und den Domains von den bekanntesten E-Mail-Providern generiert. Dazu gehören *GMX*, *WEB.DE*, *Gmail*, *T-Online*, *Freenet* und *1&1*. [Anb19]

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 dadurch die Anzahl der zu generierenden Adressen um einen großen Teil zu verringern.

Schon bei der Suche von personenbezogenen Daten wird ebenfalls nach E-Mail-Adressen gesucht. Dadurch könnte bereits eine bis jetzt unbekannte Anzahl von Adressen bekannt sein und müssten deswegen nicht mehr generiert werden.

Die erzeugten Adressen werden anschließend auf Validität geprüft. Hierfür gab es früher eine *VRFY* Anfrage von SMTP. Mit dieser Anfrage konnte eine angegebene E-Mail-Adresse überprüft werden. Allerdings wurde der Dienst von Spammern ausgenutzt und

wird dadurch von den meisten SMTP-Servern nicht mehr zu Verfügung gestellt. [BPH⁺10] Demnach muss die Validität auf einem anderen Weg geprüft werden. Eine Möglichkeit zur Prüfung ist die Verwendung bereitgestellter Webseiten, bei der die zu prüfenden E-Mail-Adresse angegeben werden kann. Eine anschließende Rückmeldung verrät dann, ob die Adresse verwendet wird oder nicht. Eine Webseite dafür wäre "<https://centralops.net/co/>". Als Alternative dazu, ist die Entwicklung eines Skriptes, welches die Validität der Adresse prüft.

Im Fall, dass mehrere Adressen von diesem Adresspool gültig sind, kann nach Einträgen diese E-Mail im Internet gesucht werden. Falls es eine Übereinstimmung mit der Zielperson gibt, wird diese E-Mail ausgewählt. Andernfalls wird an jede gültige Adresse eine Phishing-Mail gesendet. ???Kann insta,xing nach emails suchen???

5.5 Erstellung der E-Mail-Muster

Für die Erstellung der E-Mail-Muster kann eine eigene Klasse erstellt werden, welche für die Erzeugung des Textes zuständig ist. In dieser Klasse werden Strings gespeichert die einem Lückentext ähneln. Abhängig von den gefundenen Daten wird ein Lückentext gewählt, welcher anschließend mit den Daten an den passenden Lücken ergänzt wird. Mit dieser Methode muss jedoch für jede Kombination aus gewonnenen Daten ein Lückentext vorhanden sein.

Grundsätzlich können die Muster in zwei große Kategorien unterteilt werden. Es gibt einen privaten und geschäftlichen Teil. Der private Teil hat weiter Unterteilungen wie beispielsweise Familie, Hobby und Interessen. Der Text kann hier in einer Alltagssprache erstellt werden. Für ein geschäftliches Muster sollte eine gehobene Sprache angewendet werden und Daten wie der Firmenname muss bekannt sein.

5.6 Erzeugung der Phishing-Mail

Es kann je nach gewonnener Information der Person entschieden werden welches Muster gewählt werden soll.

6 Bewertung der Lösungsideen anhand der Anforderung

Mit der Programmiersprache Python lässt sich das Programm entsprechend den Anforderungen entwickeln und es kann sowohl eine Konsolenanwendung als auch eine Oberflächenanwendung programmiert werden. Es bringt alle Module mit sich um das Projekt mit dem vorgegebenen Zielen umzusetzen. Außerdem eignet sich Python sehr gut für die Bearbeitung von linguistischen Daten. [BKL09]

Um möglichst viele Informationen über eine Person im Internet zu finden, bietet die Personensuche, welche sich abhängig von den eingegebenen Daten variieren kann, die Lösung mit den meisten Vorteilen. Unter anderem kann die Arbeit des web crawlings ausgelagert werden, da nur noch die Suchergebnisse analysiert werden müssen. Allerdings muss beachtet werden, dass Benutzern bei verschiedensten Social-Media-Seiten auswählen können, ob das Benutzerprofil von einer Suchmaschine gefunden werden kann oder nicht. Aus diesem Grund, werden bei dieser Suchart die Ergebnisse kontrolliert ob sich die geforderten Seiten darin befinden. Wenn das nicht der Fall ist, wird separat auf diesen Seiten nach Information gesucht. Zu den geforderten Seiten zählen beispielsweise *XING* und *LinkedIn*.

Für die Bewertung der Lösungsideen zur Frage, wann es sich um die gesuchte Person handelt in Kapitel 5.2.2, gilt, dass alle Ideen eine Verbesserungen des Ergebnisses mit sich bringen. Allerdings gibt es Unterschied in der Wirksamkeit und in der Laufzeit des Programms. Die Erweiterung der Kriterien 5.2.2 bringt keine große Laufzeitänderung mit sich und stellt eine sehr gute Eigenschaft zur Optimierung der Informationsfindung dar, da die Zeit ebenfalls mit einbezogen wird.

Wenn die Kontakte der Suchperson in Betracht gezogen werden, kann erkannt werden wann es sich um die gesuchte Person handelt. Darüber hin Für die optimal Informationsbeschaffung einer ausgewählten Person eignet sich die Methode der Automatic Keyword Extraction um die Information wird bei der Informationsbeschaffung einer ausgewählten Person der Ansa !!XING kann man angeben ob man durch google gefunden wird!!!

Die Suchfunktion für eine große Anzahl von Personen kann *hartkodiert* werden und benötigt dadurch keine Textanalyse, da der Aufbau der Webseite im voraus bekannt ist. Das bedeutet, dass das Programm genau weiß wo welche Information auf einer Webseite steht. Auf der Seite “*www.fupa.net*“ befindet sich beispielsweise der Name einer Person immer an der gleichen Position einer Tabelle. Das bringt den Vorteil mit sich, dass der Text nicht analysiert werden muss und das Programm genau weiß, was mit diesen Daten gemacht werden muss. Zusätzlich entsteht eine sehr performante Methode zur Auslesung von personenbezogenen Daten.

Für die E-Mail-Adressgenerierung wird ein eigener Algorithmus entwickelt. Im Gegensatz zu dem Open Souce-Tool [Baz18] besteht bei diesem Algorithmus eine höhere Wahrscheinlichkeit, dass die richtige E-Mail-Adresse enthalten ist, da das Geburtsjahr, falls es bekannt ist, mit einbezogen wird. Für eine bessere Laufzeit des Programms, wird ein Skript zur Überprüfung der Adressen auf Verfügbarkeit und Gültigkeit, verwendet.

7 Umsetzung

7.1 Textanalyse mit Hilfe von Python NTLK

Mit dem *Natural Language Toolkit* ist es möglich, den vorhandenen Webseitentext zu analysieren. Zu Beginn können sogenannte “stopwords” aus dem vorgegebenen Text herausgefiltert werden. Stopwords sind Wörter die sehr oft auftreten und keinen großen Informationsgewinn mit sich bringen. Beispiele dafür sind ist, ein, einer, usw. Dadurch verringert sich die Anzahl der gesamten Wörter im Text um einen sehr großen Teil. Anschließend können Funktionen wie das Zählen des Vorkommens einzelner Wörter angewendet werden, um einen Überblick von dem Text zu bekommen. Des Weiteren kann der Text in Fragmente zerlegt werden um weitere Informationen über den Inhalt zu erlangen. Abschließend kann eine Liste der analysierten Wörter bzw. Fragmente erstellt werden. Für die Erkennung wichtiger Schlüsselwörter Es wäre denkbar, Datenbanken bzw. Wortsammlungen zu erstellen, welche die zu suchenden Schlüsselwörter beinhalten. Mit diesen Datenbanken kann nun die Liste mit den bereits verarbeiteten Wörter verglichen werden. Die Datenbanken können mit Hilfe von bekannter Listen im Internet befüllt werden. Beispiele hierfür sind eine aktuelle Liste aller Hochschulen in Deutschland, Berufsbezeichnungen, Studiengänge, Hobbys, Städte und Gemeinden, etc..

7.2 Informationsbeschaffung von der Website `www.fupa.net`

7.2.1 Erstellung eines Web Crawlers

Anforderung

Der Web Crawler soll die komplette Webseite `www.fupa.net` durchgehen und Links mit Spielerinformationen speichern. Die Funktionsweise des Web Crawlers besteht darin, dass das Programm auf der Startseite von Fupa.net beginnt nach links zu suchen und diesen folgt.

Probleme

1. Python hat einen verkürzten und erkennbaren Standard http-Header. Dieser wird von vielen Administratoren geblockt und mit der Fehlermeldung 451 erkennbar gemacht. 451 for legal reason
2. Honeypots gewollt oder ungewollt, hier Kalender darstellung mit links zu neuen Jahren die eine sehr hohe bis überhaupt keine Begrenzung haben.
3. Rekursion erreicht schnell die Maximale tiefe von 1500.
4. Zu langsamer Algorithmus

Lösungen

1. http-Header selber konfigurieren
2. Links mit möglichen Honeypots nicht beachten
3. Stack Klasse schreiben damit keine Rekursion benötigt wird
4. Algorithmus anpassen auf fupa-Webseite

7.3 Datenverwaltung und Speicherung

7.3.1 Speicherung von Personendaten in CSV oder MySQL

8 Evaluation der Implementation

9 Hauptteil

9.1 Hauptteil

9.1.1

10 Hauptteil

10.1 Hauptteil

10.1.1

11 Ethische und rechtliche Betrachtung

11.1 Ethische Betrachtung

<http://analysis.seclab.tuwien.ac.at/papers/raid2010.pdf>

11.1.1

12 Schlussbemerkungen und Ausblick

A Ein Kapitel des Anhangs

Glossar

Active Directory

Active Directory ist in einem Windows Server 2000, Windows Server 2003, oder Windows Server 2008-Netzwerk der Verzeichnisdienst, der die zentrale Organisation und Verwaltung aller Netzwerkressourcen erlaubt. Es ermöglicht den Benutzern über eine einzige zentrale Anmeldung den Zugriff auf alle Ressourcen und den Administratoren die zentral organisierte Verwaltung, transparent von der Netzwerktopologie und den eingesetzten Netzwerkprotokollen. Das dafür benötigte Betriebssystem ist entweder Windows Server 2000, Windows Server 2003, oder Windows Server 2008, welches auf dem zentralen Domänencontroller installiert wird. Dieser hält alle Daten des Active Directory vor, wie z.B. Benutzernamen und Kennwörter. 3

Glossareintrag

Erweiterte Informationen zum einem Wort oder einer Abkürzung, ähnlich einem Eintrag im Duden. 3

Abkürzungsverzeichnis

AD Active Directory 3

Symbolverzeichnis

π Die Kreiszahl. 3

Literatur

- [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.
- [Baz] BAZZELL, MICHAEL: *Email Assumptions*. <https://inteltechniques.com/osint/email.html>. Abrufdatum: 01.02.2019.
- [Baz18] BAZZELL, MICHAEL: *Open Source Intelligence Techniques: Resources for Searching and Analyzing Online Information*. CreateSpace Independent Publishing Platform, USA, 6th , 2018.
- [BKL09] BIRD, STEVEN, EWAN KLEIN EDWARD LOPER: *Natural language processing with Python: analyzing text with the natural language toolkit*. Ö'Reilly Media, Inc.“, 2009.
- [BPH⁺10] BALDUZZI, MARCO, CHRISTIAN PLATZER, THORSTEN HOLZ, ENGIN KIRDA, DAVIDE BALZAROTTI CHRISTOPHER KRUEGEL: *Abusing social networks for automated user profiling. International Workshop on Recent Advances in Intrusion Detection*, 422–441. Springer, 2010.
- [Bun18] BUNDESKRIMINALAMT: *Polizeilich erfasste Fälle von Cyberkriminalität im engeren Sinne* in Deutschland von 2004 bis 2017*. <https://de.statista.com/statistik/daten/studie/295265/umfrage/polizeilich-erfasste-faelle-von-cyberkriminalitaet-im-engeren-sinne-in-deuts> 2018. Abrufdatum: 29.10.2018.
- [Cal13] CALDWELL, TRACEY: *Spear-phishing: how to spot and mitigate the menace*. Computer Fraud & Security, 2013(1):11–16, 2013.

- [CH15] CHRISTOPHER HADNAGY, MICHELE FINCHER: *Phishing Dark Waters: The Offensive and Defensive Sides of Malicious E-mails*. 2015.
- [DSG] DSGVO: *Art. 4 DSGVO Begriffsbestimmungen*. <https://dsgvo-gesetz.de/art-4-dsgvo/>. Abrufdatum: 09.01.2019.
- [EAD09] ELDESOUKI, MOHAMED I, W ARAFA K DARWISH: *Stemming techniques of Arabic language: Comparative study from the information retrieval perspective*. The Egyptian Computer Journal, 36(1):30–49, 2009.
- [Fir] FIREEYE, INC: *Spear-Phishing-Angriffe ? Warum sie erfolgreich sind und wie sie gestoppt werden können*.
- [Had11] HADNAGY, CHRISTOPHER: *Social Engineering: The Art of Human Hacking*. 2011.
- [Jam05] JAMES, LANCE: *Phishing Exposed: Uncover Secrets from the Dark Side*. 2005.
- [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.
- [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.
- [NW18] NORDRHEIN-WESTFALEN, VERBRAUCHERZENTRALE: *Phishing-Radar: Aktuelle Warnungen*. <https://www.verbraucherzentrale.nrw/wissen/digitale-welt/phishingradar/phishingradar-aktuelle-warnungen-6059>, 2018. Abrufdatum: 29.10.2018.
- [RECC10] ROSE, STUART, DAVE ENGEL, NICK CRAMER WENDY COWLEY: *Automatic keyword extraction from individual documents*. Text Mining: Applications and Theory, 1–20, 2010.
- [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.

-
- [uDsiNe15] NETZ E.V., DATEV UND DEUTSCHLAND SICHER IM: *Verhaltensregeln zum Thema "Social Engineering"*. 2015.