



Hochschule **RheinMain**
Fachbereich Design Informatik Medien
Studiengang Medieninformatik

Abschlussarbeit

zur Erlangung des akademischen Grades

Master of Science

Data Mining komplexer Datenstrukturen aus PDF-Dokumenten

Vorgelegt von Deniz Aydar
am 25. Juli 2022
Referent Prof. Dr. Dirk Krechel
Korreferent Prof. Dr. Philipp Schaible

Erklärung gemäß ABPO

Ich erkläre hiermit, dass ich

- die vorliegende Abschlussarbeit selbstständig angefertigt,
- keine anderen als die angegebenen Quellen benutzt,
- die wörtlich oder dem Inhalt nach aus fremden Arbeiten entnommenen Stellen, bildlichen Darstellungen und dergleichen als solche genau kenntlich gemacht und
- keine unerlaubte fremde Hilfe in Anspruch genommen habe.

Wiesbaden, 25. Juli 2022

Deniz Aydar

Erklärung zur Verwendung der Masterthesis

Hiermit erkläre ich mein Einverständnis mit den im folgenden aufgeführten Verbreitungsformen dieser Abschlussarbeit:

Verbreitungsform	Ja	Nein
Einstellung der Arbeit in die Hochschulbibliothek mit Datenträger		×
Einstellung der Arbeit in die Hochschulbibliothek ohne Datenträger		×
Veröffentlichung des Titels der Arbeit im Internet	×	
Veröffentlichung der Arbeit im Internet		×

Wiesbaden, 25. Juli 2022

Deniz Aydar

Zusammenfassung

PDF-Dokumente besitzen viele Informationen, aus denen sich neue Daten generieren lassen können. Doch das Extrahieren von solchen Daten ist heutzutage immer noch mit Hürden verbunden. Dies gilt auch für die Dokumente, die in den Prozessen von Autohäusern und deren Kfz-Werkstätten verwendet und anschließend gelagert werden. Um die Frage zu beantworten, inwiefern neue Daten aus dieser Art von Dokumenten verarbeitet werden können, wird im Rahmen dieser Arbeit ein System konzipiert und entwickelt, welches es ermöglichen soll, weiteres Wissen zu gestalten beziehungsweise zu generieren. Dieses System will dabei Ansätze aus der künstlichen Intelligenz nutzen und eine grafische Schnittstelle für die Möglichkeit der menschlichen Nutzung anbieten.

Inhaltsverzeichnis

1 Einleitung	2
1.1 Zielsetzung	3
1.2 Ablauf	5
2 Hintergrund	7
2.1 Kontext	7
2.2 Analyse	8
2.2.1 Extraktionsbibliotheken	9
2.2.2 Dokumentarten	11
2.2.3 Frameworks	16
3 Konzeption und Umsetzung	17
3.1 Aufbau	17
3.2 Symbolische KI	17
3.3 Entwicklungsumgebung	17
3.4 Machine Learning	17
4 Ergebnisse	18
5 Konklusion und Ausblick	19

Kapitel 1

Einleitung

„Digitalisierung im Alltag voranbringen“ – Das war einer der Wahlslogans während der Bundestagswahl 2021. Gemeint war damit die Forderung nach einer zunehmende Digitalisierung im privaten Alltag vieler Bürger:innen, aber auch in der Wirtschaft machte sich wachsend der Wunsch nach mehr digitalen Alternativen breit (vgl. Bundesregierung 2021). Dieser Wunsch beinhaltete vor allem einen Wechsel von gängigen Papierformen verschiedenster Dokumente hin zu denselben in digitaler Ausprägung.

Genau jenes Bedürfnis nach Digitalisierung betrifft auch die vielen Arbeitnehmer:innen und Händler:innen in Autohäusern und deren Kfz-Werkstätten, die durch ihre beruflichen Tätigkeiten mit einer Vielzahl an Unterlagen, Dokumenten oder Belegen arbeiten müssen. Unter dieser Vielzahl fallen Dokumente wie Werkstatt-, Kauf- und Mietverträge sowie Rechnungen oder Diagnoseberichte.

Jene Beschäftigte in einigen Autohäusern und deren Kfz-Werkstätten sind Kund:innen bei ilexius GmbH. Ilexius bietet Enterprise-Resource-Planning (ERP) Systeme an, mit Hilfe derer Arbeitsaufträge und Arbeitsschritte innerhalb des täglichen Arbeitsprozesses in Autohäusern und deren Werkstätten durch Digitalisierung vereinfacht werden. In Kooperation von ilexius GmbH werde ich daher die Problematik meiner Thesis, die ich im folgenden noch genauer beschreiben werde, lösen und in die Tat umsetzen.

Durch eine Digitalisierung jener Dokumente eröffnen sich Vorteile wie bessere Zugänglichkeit, größere Langlebigkeit und vor allem eine angepasste und leichtere Nutzung, die zu einer höheren Effektivität innerhalb täglicher Arbeitsschritte führt. Durch bisherige erste Schritte der Digitalisierung sind diese benötigten Dokumente bereits elektronisch aufgearbeitet und den Beschäftigten in Autohäusern und deren Kfz-Werkstätten zur Verfügung gestellt worden.

Die Folge dessen ist, dass alle Dokumente standardisiert sind und dadurch die in den Dokumenten beinhalteten Informationen neu verarbeitet werden können. Eine Extraktion der Daten der einzelnen Dokumente ist jedoch nur eingeschränkt möglich, da die Struktur im gängigen Gebrauch im Dateiformat Portable Document Format (PDF) festgelegt ist.

Jene Limitierung der Datenextraktion wirkt sich gleichermaßen auf die Dealer-Management-Systeme (DMS), mit solchen die Mitarbeiter:innen ihre Prozesse abwickeln, aus. Der daraus resultierende Arbeitsaufwand, welcher dabei entsteht, um die Daten wiederverwendbar zu konstruieren, ist derzeit enorm.

1.1 Zielsetzung

Aus diesem Grund möchte ich innerhalb dieser Master-Thesis ein System entwickeln, das genau jene Problematik erleichtert und löst. Innerhalb üblicher Methoden werden heutzutage die Daten zunächst über eine grafische Anzeige mithilfe einer bereits existierenden Benutzeranwendung via gängiges Kopieren und Einfügen entnommen. Dieser Schritt funktioniert zwar im Regelfall, ist jedoch oftmals stark fehleranfällig und kann lückenhaft sein.

In so einem Fall muss das Einfügen und Kopieren manuell durchgeführt werden, was bei einer riesigen Menge an Dokumenten zu einer monotonen Arbeit für die Händler der Autohäuser führt. Andernfalls können die Daten eigenhändig abgeschrieben werden – jener Vorgang benötigt jedoch viele Ressourcen. Eher geeignet ist es, einen Datenkonverter für die PDF Dateien zu nutzen, um so die Dokumente beispielsweise in Bilder umzuwandeln, wodurch die Daten zugänglicher sind, aber immer noch verarbeitet werden müssen. Die letztgenannte Möglichkeit wirkt allerdings eher wie eine Übergangslösung als eine professionelle endgültige Durchführung.

Aufgrund dieser bisherigen teils aufwendigen und mit Fehlern verbundenen Möglichkeiten möchte ich mich in meiner Thesis von diesen Optionen abwenden und das Data-Mining nutzen. Das Data-Mining soll eine Automatisierung unterstützen, mit der Daten aus dem PDF-Dokument extrahiert werden. Allein über das Data-Mining ist es möglich, Inhalte zu extrahieren und für andere Systeme bereit zu stellen.

Der Unterschied zwischen meiner Verwendung des Data-Minings und zum gängigen Data-Mining liegt dabei in den Einschränkungen des Dateiformats: der Quelltext einer solchen Datei stellt erstens keine klare Hierarchie der Daten dar und zweitens besitzt es durch das Fehlen von Markierungen keine Informationen darüber, was die Daten an sich darstellen sollen [11].

Mit Hilfe dieser Technik des Data-Minings soll es ermöglicht werden, aktuelle Prozesse detaillierter zu steuern und zu überwachen. Des Weiteren können unter anderem Abgleiche von Rechnungen mit dem System durchgeführt werden und weitere Datenanalysen und Reporting kreiert werden. Mit dem derzeitigen Stand der Datenextraktion können jedoch lediglich Agierende aus dem technischen Bereich arbeiten, da allein sie das dafür technisch notwendige Know-how besitzen. Die Technologien hierfür benutzen unterschiedliche Ansätze und müssen daher zunächst für diesen Fall ausgewählt werden.

Damit jedoch auch für Agierende innerhalb des technischen Bereichs die Nutzung nicht zu abstrakt bleibt, soll durch die Entwicklung eines Systems der Zugriff greifbarer gestaltet werden, welches wiederum allein durch eine benutzerfreundliche Anwendung ermöglicht wird. Die Benutzer sollen eine Möglichkeit bekommen und über eine Eingabe der Steuerung bestimmten können, welche Informationen extrahiert werden sollen.

Darüber hinaus muss eine Automatisierung vorhanden sein, um auch eine große Datenmenge verarbeiten zu können; zudem sollen die Vorgänge eine niedrigere Fehlerquote aufzeigen. Die Prozesse müssen außerdem während des gesamten Ablaufs über eine Steuerung und Regelung kontrollierbar sein. Solche Prozesse können zum Beispiel über Regeln festgelegt beziehungsweise gesteuert werden, die bei einer Erfüllung weitere Aktionen oder Regeln auslösen können.

Diese beschriebene Eigenschaft kann durch eine symbolische künstlichen Intelligenz (KI) abgedeckt werden, welche eine vorgegebene Verarbeitung möglich macht. Hiermit bezeichnet man einen altmodischen Ansatz der KI, bei der über das Festlegen von Symbolen menschliches Wissen in einer Logik gehalten wird und diese benutzt wird um weiteres Wissen zu generieren [10]. Durch eine Parametrisierung einer solchen KI kann dann die Unschärfe der ausgewählten Daten festgelegt werden, damit eine Feineinstellung stattfinden kann. Das heißt, der Benutzer eines solchen Systems soll sich beginnend von groben Definitionen für die Verarbeitung zu einer detaillierten und angepassteren Definition über Feedback des Systems hinarbeiten.

Mit Hilfe dieser Annäherung an Definitionen und Regeln kann so für eine bessere Erfolgsquote gesorgt werden. Zudem kann der Ansatz des Machine Learnings (ML), wobei erfolgreiche Verarbeitungen einem System antrainiert werden, weitere Dokumente aus Daten extrahieren. Hierbei ist noch offen, welcher Typ des Machine Learnings sich für diesen Fall eignet.

Beide Ansätze – der der symbolischen KI und der des Machine Learnings – bieten als Option die Entwicklung einer grafischen interaktiven Entwicklungsumgebung (IDE) an. Jene soll das Festlegen von Definitionen und Regeln erlauben,

mit welchen die Dokumente verarbeitet werden können. Außerdem soll diese Oberfläche verschiedene Funktionalitäten anbieten und auch Feedback sowohl bei einem problemlosen Lauf als auch bei einem fehlerhaften Lauf zurückgeben.

Des Weiteren soll das System es ermöglichen Änderungen der Definitionen anzumerken und Unterschiede zu erstellen, womit auch bisherige Ergebnisse angezeigt werden. Die Festlegung von Definitionen und Regeln soll außerdem durch eine Ansicht der Dokumente unterstützt werden. Insgesamt wird das System die zwei Ansätze als Subsysteme aufteilen um die Umsetzung zu modularisieren und einen Vergleich zu ermöglichen.

1.2 Ablauf

Um genau dieses optimale System für das beschriebene Szenario entwickeln zu können, werde ich in dieser Arbeit wie folgt vorgehen:

Für eine vollständige Auseinandersetzung soll eine weitere Analyse stattfinden, damit die vorhandenen Grundlagen für den Anwendungsfall evaluiert werden können. Das heißt, es werden mit Hilfe der Funktionalitäten der Bibliotheken Ergebnisse auf Basis der Datensätze erzeugt, um die Erfolgsquoten zu überprüfen. Diese werden dann verglichen und ausgewertet. Parallel dazu soll die Umgebung für die Entwicklung eingerichtet werden, mit der die Implementierung der Systeme stattfinden soll.

Danach widme ich mich der Konzeption und der Entwicklung des Systems. Diese beginnt mit der symbolischen KI als ein Subsystem des gesamten Projekts, sodass ein direkterer Ansatz ermöglicht wird. Die Entwicklung hierbei wird im Wasserfall-Ansatz, einem Ansatz bei der phasenweise die Eigenschaften einer solchen KI umgesetzt werden, um eine Grundlage für die nächste Komponente zur Verfügung zu stellen.

Bei dieser Komponente handelt es sich um die interaktive Entwicklungsumgebung, die den Zugang für den Benutzer zum System darstellt. Deswegen folgt im nächsten Schritt die Konzeptionierung und Entwicklung der interaktiven Entwicklungsumgebung, welche sich während der Entwicklung der symbolischen KI parallelisieren lässt. Sobald das Gerüst der Benutzeroberfläche fertiggestellt ist, möchte ich dies mit der symbolischen KI verbinden.

Als Gegenstück wird sich dann mit dem Machine Learning Subsystem, bei der gleichermaßen eine Konzeptionierung und Entwicklung stattfindet, gewidmet. Hier wird das bisherige Gesamtkonstrukt in einem iterativen Ansatz umgesetzt, sodass der Hauptteil so früh wie möglich verfügbar ist, damit das Training

des ML-Systems durch die Dokumente auch zeitnah starten kann. Auf das Auswerten dieses Subsystems folgt eine Anpassung dessen.

Danach wird das Software Testing für die beiden Subsysteme eingeführt. Mit Unit Testing, Integration Testing und Functional Testing wird das erfolgreiche Ausführen gewisser Eigenschaften abgedeckt. Dies dient wiederum als Grundlage dafür, die Continuous Integration (CI) für das System einzuführen, womit eine sichere und konsistente Entwicklung nach dem Prototyping angeboten wird.

Wenn dann aber die Systeme gegeben sind, kann das grafische Tool weiter angepasst werden, welches sich dann mit den Systemen verbinden soll, um die Ausführbarkeit der Funktionalitäten zu überprüfen. Um ein gemeinsames System zu ermöglichen, welches dann gegebenenfalls weitere Testschritte bedarf, werden dann die Komponenten zusammengeführt. Sobald jenes eine erfolgreiche Form annimmt, wird der Teil der Continuous Integration mit einbezogen, mit jener das System einen Zustand der Instandhaltung annehmen kann. Über Testing sollen so mögliche Fehlerquellen entdeckt werden, bevor diese überhaupt auftreten können.

Unter die Zielsetzung für meine beiden eben genannten Ansätze fallen die Kundendaten mit den Eigenschaften wie Name, Firma, Adresse und Kundennummer, die Fahrzeugdaten worunter Fahrgestellnummer, Modell, Farbe, Motor, Erstzulassung, Letzter/nächster Service und TÜV fallen und die Abrechnungsdaten wie die Jobs mit Arbeitspositionen und Teilen, Preise etc.

Dazu stehen als weitere Grundlage mehr als 10.000 annotierte unterschiedliche Dokumente mit genau den Daten, die extrahiert werden sollen, als Datenbank zur Verfügung. In Zukunft sollen für das Mengengerüst mehr als 100.000 Dokumente pro Jahr bearbeitet werden. Hierbei ist es aber auch möglich, aus den bereits bestehenden Daten einen Pool für den Lernprozess bzw. für die Verifikation zu bilden.

Zum Schluss werde ich die Ergebnisse des gesamten Systems betrachten und einen Ausblick zu der Thematik geben.

Kapitel 2

Hintergrund

Um die Problematik detaillierter darzustellen und um eine Übersicht zu ermöglichen, wird in diesem Kapitel der Rahmen der Problematik aus technischer Sicht genauer dargestellt. Hierbei werden die möglichen Technologien abgewogen und evaluiert. Zunächst jedoch wird der Kontext ausführlicher erklärt.

2.1 Kontext

Im Automobilbereich wird der Verkauf von Kraftfahrzeugen in zwei Segmente aufgeteilt: dem Sales-Bereich, bei dem es sich um den Verkauf von Neu- und Gebrauchtwagen handelt und dem After-Sales-Bereich, bei jenem eine Bindung zum Kunden über den Verkauf von Verschleiß- und Ersatzteilen geschaffen wird. [17] Diese Aufteilung der Segmente wird ebenfalls von den Autohäusern angewendet, zugleich werden für diese jeweils gängige Abläufe als Prozesse benutzt.

Für den After-Sales-Bereich bietet ilexius GmbH zwei Arten von ERP-Systeme an. Die erste Art ist das von der Automobilmarke Jaguar Land Rover finanzierte ‘vTab’, welches eine Plattform für elektronische Fahrzeugkontrollen und ein Backend für Jaguar Land Rover Mitarbeiter:innen, mit dem die Autohäuser bewertet werden können, anbietet. Die zweite Art mit dem Namen ‘ATT’ ist eine Plattform. Mit derjenigen können Autohäuser ihre Arbeitsprozesse in Schrift, Sprache und Bild dokumentieren. Die Dokumentation erfolgt mittels mobiler Endgeräte und werden an das ERP-System weitergeleitet, sodass die Dokumentation archiviert werden kann.

Die ERP-Systeme, die bereits im Kapitel 1.1 beschrieben worden sind, unterstützen damit die Kunden:innen von ilexius GmbH in den Autohäusern und

deren Werkstätten mit der Planung, der Steuerung sowie der Verwaltung von verschiedenen Aufgaben in ebenjenem After-Sales-Bereich.

Eine Art der Aufgaben sind Arbeitsaufträge, bei der es sich in den meisten Fällen um die Behandlung eines Autos inklusive Kundendaten handelt. Der Auftrag besteht hierbei aus einer Liste an Arbeitsabläufen, die als Jobs bezeichnet werden. Die Jobs wiederum unterteilen sich in einzelne Arbeitsschritte, die Joblines genannt werden, welche die konkreten Verrichtungen der Arbeit darstellen und somit den gesamten Arbeitsauftrag in einzelne Teilschritte vervollständigen.

Sind diese Aufträge abgeschlossen, benutzen die Mitarbeiter:innen der Autohäuser ihre Dealer-Management-Systeme, um eine Rechnung zu erstellen und diese im ERP-System zu archivieren. Das ERP bietet des Weiteren eine Schnittstelle für diese Dokumente, wie z.B die Rechnungen oder die Garantieaufträge an, die allein über Scanner die Dokumente übertragen. Dort werden die Dokumente verarbeitet und mit Hilfe eines OCR Scanners digitalisiert.

Aktuell ist allerdings die Digitalisierung nicht in der Lage, einen Kontext beziehungsweise (bzw.) eine Semantik für diese Dokumente zu erstellen. Idealerweise kann dies genutzt werden, um aus den Dokumenten Arbeitsaufträge zu generieren und somit zu digitalisieren. Demzufolge ergibt sich hier an der Stelle die Möglichkeit, einen neuen Ansatz zu testen, sodass die Informationen aus diesen Dokumenten als Daten erhaltbar sind.

2.2 Analyse

Das Extrahieren von Daten ist eine gängige Methode um Informationen aus verschiedenen Systemen wie Datenbanken oder Software as a Service (SaaS) Plattformen zu beschaffen [15]. Durch diese Beschaffung können die Daten weiter verarbeitet werden um neue Lösungsmöglichkeiten für das eigene System anzubieten. Die Arten des Extrahierens unterscheiden sich hierbei im Zeitverlauf und der Herangehensweise. So können Diese extrahiert werden sobald Änderungen der Daten beobachtet oder mitgeteilt werden, wodurch die Daten nach und nach beschafft werden.

Des Weiteren gibt es die Option eine komplette Extraktion durch zu führen. Dieser Ansatz kann sich in vielen Fällen als nachteilig erweisen, da die Daten jedes mal bei Anpassungen der ursprünglichen Information neu verarbeitet werden müssen. Da sich die Problematik im Rahmen der Thesis auf Dokumente, die bereits im ERP-System archiviert wurden, bezieht, trifft dieser Fall hier nicht zu.

2.2.1 Extraktionsbibliotheken

Für das Data Mining von PDF Dateien existieren bereits Lösungen, die von großen Unternehmen wie Amazon Web Services (AWS) oder Adobe bereits angeboten werden. So nutzen Produkte wie Textractor [5] von AWS oder Adobes Extractor API [9] Künstliche Intelligenz, Machine Learning und Optical Character Recognition (OCR) um die Extraktion der Daten zu ermöglichen. Unternehmen wie ABBYY [4] habe sich bereits auf diesem Gebiet spezialisiert und sind erfolgreich mit der Datenextraktion. Auch existieren bereits Open-Source Bibliotheken, die sich mit dem Entnehmen der Daten aus PDF Dokumenten beschäftigen.

Die Meisten davon lassen sich hinsichtlich der Aufteilung der Aufgaben gruppieren. So erlauben unter anderem PDFMiner [18] oder Apache Tika [6] Texte aus PDF Dokumenten zu gewinnen. Funktionalitäten für das Extrahieren von Daten aus Tabellen bieten Projekte wie Tabula [16] oder Camelot [8] an womit ein weiterer Aspekt durch die Nutzung abgedeckt wäre.

Bei Beiden zeigen sich niedrige Fehlerquoten auf [7] und die Nutzung von einer Fuzzy Logik. Camelot ermöglicht dazu noch ein Web Interface womit die simple Extraktion über einen Browser lokal stattfinden kann. Andere Bibliotheken wie Textricator [12] sind umfangreicher gestaltet und ergeben geeignete Schnittstellen für die Anfragen spezifischer Inhalte.

Doch auch wenn dies mit textbasierten Dokumenten funktioniert und der Export flexibel ist, können zunächst weder nicht-textbasierte noch komplexere PDF Dokumente verarbeitet werden. Des Weiteren sind die Verarbeitungen nicht automatisierbar und die Nutzergruppe wird durch eine fehlende grafische Komponente eingeschränkt. Aus dieser Problematik kommend, zeigt sich, dass community-angetriebene Projekte wie PdfMiner.six [13], welches eine Abspaltung des bereits erwähnten PdfMiner ist, und PdfPlumber [14], welches wiederum auf PdfMiner.six basiert, als geeignete Alternativen darstellen.

PdfPlumber stellt im Vergleich zu seinem Ausgangsprojekt dazu noch eine Tabellenextraktion bereit und bietet für beide Aspekte der Extraktion eine Schnittstelle an. Es hält eine Kommandozeile, die Möglichkeit, visuelles Debugging zu nutzen, und Funktionalitäten wie das Entfernen von Duplikaten in Texten, bereit.

Das visuelle Debugging, welches in Abbildung 2.1 zu sehen ist, ermöglicht es unter anderem, die Texte des Dokuments hervorzuheben, was wiederum die Extraktion greifbarer darstellt und die Entwicklung mit den Funktionalitäten unterstützt.

In dieser Abbildung sind die persönliche Daten des Auftrags verpixelt und damit anonymisiert. Auch die Funktionalität des Entfernen von Duplikaten kann dann zum Einsatz kommen, wenn Dokumente bei ihrer Erstellung oder Verarbeitung versehentlich die selben Einträge nochmal im Quellcode des PDFs besitzen.

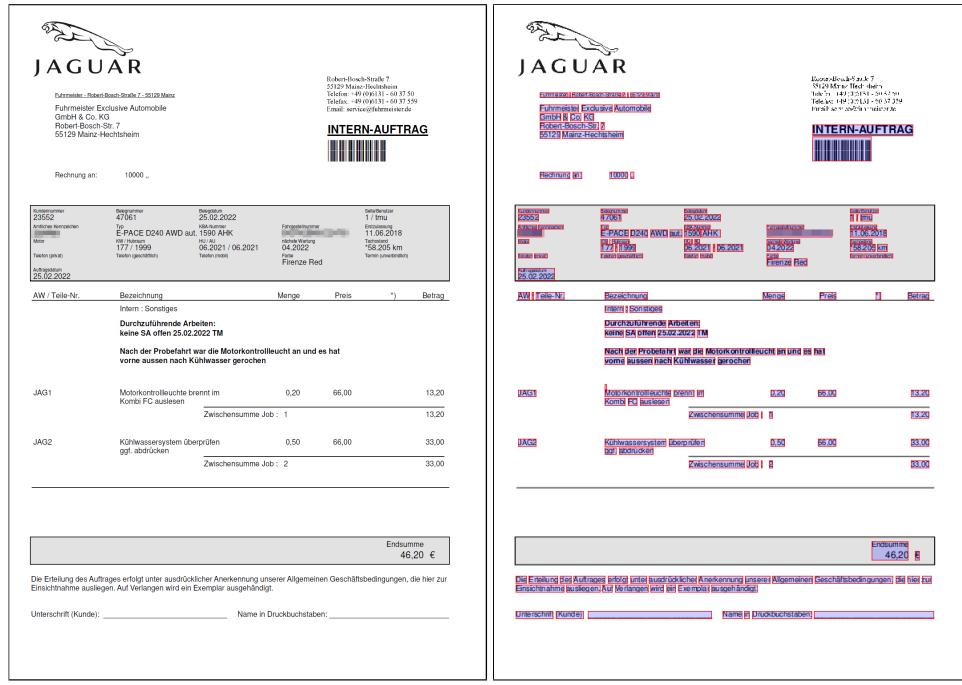


Abbildung 2.1: visuelles Debugging bei einem Dokument

Die Hürden für eine komplett automatische Extraktion bestehen trotzdem noch. So ist nämlich nicht klar, wie das zu erwartende Verhalten der Logik der Extraktion für den generellen Fall auszusehen hat. Eine differenzierte Extraktion für weitere Eigenschaften eines Dokuments wie die Zeilenumbrüche von Absätzen, die Seitennummern, die Kopf- und Fußzeilen, die Formatierung et cetera (etc.) ist daher nötig. Die Frage, die sich dabei stellt ist, welche Eigenschaften eines Dokuments genau zusätzlich extrahiert werden sollen. Um eine Auswahl zu treffen, müssen die Arten der Dokumente, die in dem ERP-System auftreten, genauer durchleuchtet werden.

2.2.2 Dokumentarten

Bei den Dokumenten, die in das ERP-System eingescannt werden, handelt es sich in der Mehrheit um Garantie- und Werkstattaufträge sowie Rechnungen. Jedes Autohaus benutzt hierbei ein individuelles Muster bzw. Format für die Arten der Dokumente. Jedoch ist der Aufbau dieser Dokumente zu Teilen indirekt durch die von den Autohaus genutzten Dealer-Management-Systeme vorgegeben.

Dadurch passiert es, dass es Autohäuser gibt, welche das gleiche DMS verwenden und sowohl das selbe als auch ein in Detail angepassten Dokumentenaufbau für ihre Aufträge benutzen. Dem Umfang dieser Thesis geschuldet werde ich mich daher auf drei Dokumententypen von unterschiedlichen Autohäusern mit ihren jeweils genutzten DMS beschränken und anhand jener drei Typen mein Vorgehen exemplarisch demonstrieren.

In Abbildung 2.2 wird die erste Art eines Auftrags dargestellt. Bei diesem digitalen Schriftstück handelt es sich zunächst von der Form her um eine Garantierechnung des Care Warranty DMS. Das Dokument besitzt hierbei drei Bereiche, die relevante Informationen für einen Arbeitsauftrag enthalten. In der oberen Hälfte befindet sich zunächst das Adressfeld, welches die Kundeninformationen im genormten Adressformat enthält, und ein weiterer Abschnitt, welches die Fahrzeugdaten festhält.

Hierbei ist zu beobachten, dass dieser Abschnitt sich mit einer Tabelle vergleichen lässt, da die einzelnen Einträge mit jeweils unterschiedlich großen Abständen verteilt dargestellt werden aber gleichzeitig noch eine Gesamthöhe pro Spalte beziehungsweise eine Gesamtbreite pro Höhe einer Tabelle erfüllen. Jedoch fällt auf, dass sich, im Gegensatz zur Anzahl der Zeileneinträge, die Anzahl der Spalteneinträge pro Zeile unterscheidet.

Außerdem fehlen jegliche Einrahmungen beziehungsweise jegliche Trennlinien für die Zeilen und die Spalten. Bei den Einträgen in diesem Abschnitt handelt es sich um die Abrechnungsnummer, das Datum der Reparatur, den Fahrzeugmodellnamen, den zuletzt abgelesenen Kilometerzählerstand, der Arbeitsauftragsnummer, der Fahrzeug-Identifikationsnummer beziehungsweise der 'Vehicle Identification Number' (VIN) und das Datum der Erstellung des Arbeitsauftrags.

Dagegen wird im unteren Bereich eine Tabelle durch die Aufzählung der einzelnen Jobs und Joblines dargestellt. Die Tabelle ist an der Kopfzeile zu erkennen, welche die Jobs und Joblines in die Arbeitsnummer beziehungsweise in die Teilenummer, in der Bezeichnung, in die Menge des Arbeitszeitwertes, in den

#!GRHMFP!#		10	70	DE
Jaguar Care Garantie				
Fahrzeughalter: Kunden-Nr.: 151365				
				
RECHNUNG GWL				
F-Type		17649	GJC	47822300
		18.03.22		
AW/T-Nr. Bezeichnung AW/Menge Preis Betrag				
00010 JAGUAR CARE WARTUNG CCC WARTUNG DURCHFÜHREN				
102111 52.000 KM (32.000 MEILLEN) ROUTINEWART	1,10	165,00	181,50	
ZZZ001 PREMIUM SERVICE PAUSCHEL			40,00	
KBEND	-----			
TERMINIERT (X) JA: NEIN:				
DATUM AUFTRAGSANNAHME:				
DATUM REPARATURBEGINN:				
DATUM REPARATUR ENDE:				
DATUM FAHRZG. ABGEHOLT:				
FRE ESRVPLNOIL	MOTORÖL	7	25,00	175,00
JAG C2D3670	OLPFILTER	1	33,92	33,92
JAG C2P2410	AIR FILTER	1	99,67	99,67
JAG C2A102010	SCHIEIBENWASCHFLUESSIGK	1	7,13	7,13
JAG AJ813202	ÖLABLA SCHRAUBE	1	5,50	5,50
008325		18.03.22	9:55	10:65
			1:10	
Summe Arbeit: 221,50				
Summe Fremdl:				
Summe Teile : 321,22				
542,72	0,00%		EUR	542,72

Abbildung 2.2: Exemplar eines Garantieauftrags

Preis und in den Betrag übersichtlich aufteilt. Der Arbeitszeitwert ist dabei eine Einheit, die für eine genaue Abrechnung der Arbeitszeiten, um eine Aufteilung der Stunden zu ermöglichen. [3] Die Menge wird hierbei als Reparaturzeit oder Anzahl angegeben. Außerdem existiert bei dieser Tabelle nur für die Kopfzeile eine untere Trennlinie. Auch hier sind Probleme bezüglich (bzgl.) der Form der Tabelle zu erkennen.

Während die ersten vier Zeilen noch unauffällig erscheinen, besitzt die fünfte Zeile bei der Bezeichnung einen längeren Text, sodass der restliche Text in weiteren Zeilen in der Spalte der Bezeichnung darunter geschrieben steht.

Während die Abstände zwischen den Spalten insgesamt konsistenter erscheinen, sind in der Abbildung 2.2 in mehreren Zeileneinträgen Überläufe bei den

Arbeitnummern bzw. den Teilenummern zu erkennen. In der letzte Zeile ist außerdem zu erkennen, dass die Einträge zu Teilen auch verschoben dargestellt werden, was eine eindeutige Erkennung der Positionierung noch weiter erschwert. Inhaltlich werden Jobs mit den dazu gehörigen Joblines darunter folgend dargestellt.

Dies ist durch die Codierung der Arbeitsnummern bzw. der Teilenummer zu erkennen und die Hierarchie ist durch das jeweilige Dealer-Management-System vorgegeben. In der ersten Spalte können daher Codes von Jobs, Reparaturzeiten, Ersatzteilen, Fremdleistungen, Stempelzeiten und Kundenbeschwerden auftreten.

Für die Codierung ist keine standardisierte Spezifikation vorhanden, weswegen dies für jedes Autohaus und deren Nutzung der DMS individuell definiert. So sind die Codes der Jobs als eine fünfstellige Zeichenkette aus Ziffern definiert, welche mit den Ziffern '00' anfangen müssen und die Reparaturzeiten besitzen eine sechsstellige Zeichenkette aus Ziffern. Unterhalb der Tabelle werden am Ende des Dokuments noch die gesamten Kosten summiert aus den Betragseinträgen dargestellt, welche wiederum sich in den Kosten der Arbeit und der Teile aufteilen lassen.

In Abbildung 2.3 wird eine Rechnung der Formel 1 DMS dargestellt. Hier ist im Vergleich zur vorherigen Abbildung 2.2 ein ähnlicher Aufbau ersichtlich. Während das Adressfeld gewöhnlich erscheint, sind jedoch die Fahrzeug- und Kundendaten diesmal in zwei unterschiedlichen Tabellen aufgeteilt.

Die Erste von diesen Tabellen befindet sich im oberen Teil der rechten Seite des Dokuments, in welchen die einzelnen Kundeninformationen aufgelistet sind. Die Tabelle für die Fahrzeuginformationen befindet sich im unteren Bereich der oberen Hälfte des Dokuments und erweist sich durch das Vorhandensein der Trennlinien als eine klassische Tabelle. Im letzten Teil der Formel 1 Rechnung befindet sich die Liste der Job und Joblines.

Hier fallen ähnliche Problematiken, wie die unterschiedlichen Abstände zwischen den Spalten der Einträge, im Vergleich zu dem Care Garantieauftrag auf. Auch existieren Überläufe der Einträge, allerdings treten diese nur für die Felder der Bezeichnung auf. Die Gesamtkosten werden dann am Ende des Dokuments in den Teilsummen der Beträge angezeigt.

Die letzten Dokumentart wird in der Abbildung 2.4 dargestellt. Hier handelt es sich um einen Werkstattauftrag von dem Dealer-Management-System KFZ3000. Dieser Auftrag lässt sich von der Aufteilung der Segmente mit dem Care Warranty Dokument geeigneter vergleichen. Während auch das Adressfeld sich wieder normiert darstellt, werden Kunden- und Fahrzeugdaten über Schlüssel-Wert

BMW Vertragshändler		ahg Autohandelsgesellschaft	BMW Service		
ahg Autohandelsgesellschaft mbH - Senefelderstr. 2 - D-73760 Ostfildern					
		Kunden-Nr. : 444515-33 Rechnungs-Nr. : 133 - Rechnungsdatum : Annemer : Auftrag/LS-Nr. : 115025 Auftragsdatum : 24.02.2022 Hauptauftrag-Nr: 112985 Seite : 1 Leistungsdatum : 24.02.2022			
TESTRECHNUNG					
Filiale : E Ostfildern					
Fahrz.Nr.	Modell/Typ	Fahrg.-Nr./EZ/HU	km Ein/Aus angen./ausgef.		
[REDACTED]	VU31 320d Touring	[REDACTED] 20.09.2006 / 09/23	177393 02321 [REDACTED] 177584 02019		
Nummer	Bezeichnung	Menge/AW/St	E-Preis	Gesamt	
0000556	Fahrzeugtest durchführen	2	10,66	21,32	
1223505	Alle Glühlampen ersetzen	16	13,63	218,08	
6100006	Prüfplan für Vorglühlammlage abgearbeitet	1	13,63	13,63	
6121528	nicht i.O. Bordnetzspannung unterstützen Bordnetzbatterie nachladen	1	13,63	13,63	
11612246945	Profilabdichtung	4	4,27	17,08	
11617790198	Profilabdichtung	4	4,27	17,08	
12237786869	Glühstift	4	14,94	59,76	
Zwischensumme				360,58	
0,00% MwSt von		360,58 =	0,00		
			Gesamt Lohn EUR 266,66 Gesamt Teile EUR 93,92		
			Endesumme EUR 360,58		
Zahlbar sofort!					
Hausanschrift ahg Entenmann - Eine Niederlassung der ahg Autohandelsgesellschaft mbH Senefelderstraße 2 73760 Ostfildern		Kontakt Telefon: +49 (0) 711/44983-0 Fax: +49 (0) 711/44983-83 E-Mail: ostfildern@ahg-entenmann.de www.ahg-entenmann.de	Bankverbindung Kreissparkasse Freudenstadt IBAN DE06 6425 1060 0000 5435 43 BIC SOLADES1FDS	Geschäftsführer Alexander Kramer Thomas Lindrich	Registergericht Stuttgart HRB 440 304 Sitz Horb a. N. USt-ID-Nr. DE811162775

Abbildung 2.3: Exemplar einer Rechnung

Anordnungen in einem größeren Feld dargestellt.

Die untere Hälfte des Werkstattauftrags besitzt auch wie die bereits anderen genannten Dokumente eine tabellenähnliche Auflistung der Job und Joblines. Auch hier sind wieder Überläufe und unterschiedliche Abstände in dem Bereich zu erkennen. Dazu kommt, dass, im Gegensatz zu den anderen Dokumentarten, zwischen den Zeilen ein Art Summenstrich hinzugefügt wird, der die Anzahl der Jobs anzeigen soll. Außerdem wird im unteren Bereich des Dokuments nur die gesamte Betragssumme dargestellt.

Insgesamt lässt sich feststellen, dass alle vorgestellten Dokumentarten mehrere Hürden als gemeinsame Problematik verbindet. Die dargebotene Darstellung

																														
<p>Fahrmeister - Robert-Bosch-Straße 7 - 55129 Mainz Fahrmeister Exclusive Automobile GmbH & Co. KG Robert-Bosch-Str. 7 55129 Mainz-Hechtsheim</p> <p>Robert-Bosch-Straße 7 55129 Mainz-Hechtsheim Telefon: +49 (0)6131 - 60 37 50 Telefax: +49 (0)6131 - 60 37 559 Email: service@fahrmeister.de</p>																														
INTERN-AUFRAG																														
																														
<p>Rechnung an: 10000 ..</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <tr> <td>Kundennummer 23552</td> <td>Beginnmar 47061</td> <td>Beigefügt 25.02.2022</td> <td>FahrerInnernummer [REDACTED]</td> <td>Seller/Benutzer 1 / Immo</td> </tr> <tr> <td>Amtliches Kennzeichen [REDACTED]</td> <td>Type E-PACE D240 AWD aut.</td> <td>Hilfs-Nr. 1590 AHK</td> <td>nächste Wartung 04.2022</td> <td>Einstellung 11.06.2018</td> </tr> <tr> <td>Motor [REDACTED]</td> <td>KW / Hubraum 177 / 1999</td> <td>HU / AU</td> <td>Farbe Firenze Red</td> <td>Tachostand *58.205 km</td> </tr> <tr> <td>Telefon (privat)</td> <td>Telefon (geschäftlich)</td> <td>Telefon (mobil)</td> <td></td> <td>Termin (unverbindlich)</td> </tr> <tr> <td colspan="5">Auftragsdatum 25.02.2022</td> </tr> </table>						Kundennummer 23552	Beginnmar 47061	Beigefügt 25.02.2022	FahrerInnernummer [REDACTED]	Seller/Benutzer 1 / Immo	Amtliches Kennzeichen [REDACTED]	Type E-PACE D240 AWD aut.	Hilfs-Nr. 1590 AHK	nächste Wartung 04.2022	Einstellung 11.06.2018	Motor [REDACTED]	KW / Hubraum 177 / 1999	HU / AU	Farbe Firenze Red	Tachostand *58.205 km	Telefon (privat)	Telefon (geschäftlich)	Telefon (mobil)		Termin (unverbindlich)	Auftragsdatum 25.02.2022				
Kundennummer 23552	Beginnmar 47061	Beigefügt 25.02.2022	FahrerInnernummer [REDACTED]	Seller/Benutzer 1 / Immo																										
Amtliches Kennzeichen [REDACTED]	Type E-PACE D240 AWD aut.	Hilfs-Nr. 1590 AHK	nächste Wartung 04.2022	Einstellung 11.06.2018																										
Motor [REDACTED]	KW / Hubraum 177 / 1999	HU / AU	Farbe Firenze Red	Tachostand *58.205 km																										
Telefon (privat)	Telefon (geschäftlich)	Telefon (mobil)		Termin (unverbindlich)																										
Auftragsdatum 25.02.2022																														
AW / Teile-Nr.	Bezeichnung	Menge	Preis	*)	Betrag																									
JAG1	Motorkontrollleuchte brennt im Kombi FC auslesen	0,20	66,00		13,20																									
	Zwischensumme Job : 1				13,20																									
JAG2	Kühlwassersystem überprüfen ggf. abdrücken	0,50	66,00		33,00																									
	Zwischensumme Job : 2				33,00																									
<div style="border: 1px solid black; padding: 5px; display: inline-block;">Endsumme 46,20 €</div>																														
<small>Die Erteilung des Auftrages erfolgt unter ausdrücklicher Anerkennung unserer Allgemeinen Geschäftsbedingungen, die hier zur Einsichtnahme ausliegen. Auf Verlangen wird ein Exemplar ausgehändigt.</small>																														
<small>Unterschrift (Kunde): _____ Name in Druckbuchstaben: _____</small>																														

Abbildung 2.4: Exemplar eines Werkstattauftrags

der Kunden- und Fahrzeuginformationen wird entweder in einer Tabellenform gemeinsam oder in zwei verschiedenen Tabellen abgebildet. Diese Tabellenform besitzt außerdem selten Trennlinien, welche geeignet sind, um eine derartige Tabelle als solche definieren zu können. Für die Job und Joblines tritt das selbe Problem auf.

Diese Auflistung wird nur durch ihre Anordnung der Einträge und deren Abstände dazwischen als solche erkenntlich, bricht in den prototypischen Dokumenten aber mindestens einmal die Form durch Überläufe oder Verschiebungen der Texte. Es ist festzustellen, dass unterschiedliche Tabellenstrukturen zu unterschiedlichen Dieses Angelegenheit besteht, da es für Tabellen keine universelle Standardisierung existiert.

Damit diese Formen durch eine Logik des Extrahierens abgedeckt werden, scheint die Möglichkeit der Fallunterscheidung zu bestehen. Bei dieser Unterscheidung sollen die Dokumente anhand ihrer Arten konzipiert und definiert werden, sodass die Entnahme detaillierter stattfinden kann. Für andere Bereiche eines Dokuments wie dem Adressfeld werden gewöhnliche Funktionalitäten der Extraktion genutzt und zur Verfügung gestellt. Diese können dann von allen weiteren Dokumentarten geteilt und genutzt werden, womit eine Verallgemeinerung verschiedener Dokumententypen ermöglicht wird.

2.2.3 Frameworks

Damit das System für den Prozess der Extraktion als solches sich zeitnah im Rahmen einer Thesis umsetzen lässt, können Frameworks dem ganzen Zeitfenster entgegenkommen. Hierbei werden sich für die Teilsysteme verschiedene Technologien angeschaut und in Erwägung gezogen.

So wird zunächst die symbolische künstliche Intelligenz dem ersten Teilsystem zugeordnet und in seiner Art unterschieden. Zudem müssen die Realisierungsmöglichkeiten der symbolischen KI betrachtet und aus den verschiedenen Arten eines wissensbasiertes Systems ein Modell ausgewählt werden. Ein solches wissenbasiertes System ist hierbei der Oberbegriff für Programme, die mit Hilfe einer Wissensbasis Mechanismen der Schlussfolgerung nutzen, um Lösungen für Probleme zu finden [2]. Es existieren verschiedene Modelle eines wissenbasierten Systems, die sich, je nach Anwendungsfall, als passend erweisen. So gibt es den Ansatz des regelbasierten Systems, bei dem die Schlussfolgerung über Regeln und Fakten, die von einem Benutzer definiert werden, gestaltet wird. Bekannt für die Erstellung eines solchen Systems ist CLIPS[1], welches über eine eigene Syntax erlaubt,

Kapitel 3

Konzeption und Umsetzung

Für den Sprung zur Umsetzung bedarf es als Nächstes die Konzeption des technischen Rahmens. Danach wird der Verlauf der Entwicklung dargestellt und die

3.1 Aufbau

3.2 Symbolische KI

Wie im Segment

3.3 Entwicklungsumgebung

Für die graphische Schnittstelle der Benutzer des Systems

3.4 Machine Learning

Während

Kapitel 4

Ergebnisse

Kapitel 5

Konklusion und Ausblick

Im Rahmen der Abschlussarbeit wurde ein System konzipiert und umgesetzt, welches eine Extraktion von Dokumenten über zwei Ansätze ermöglicht. Über die Teilsysteme wurden technische Funktionalitäten eingeführt werden um dem Kapitän die Webanwendung als PWA anzubieten.

Insgesamt eignen sich Webanwendungen um Ressourcen übersichtlicher darzustellen und mobil zu beobachten. Auch können Offline-Zustände dem Benutzer trotzdem erlauben Aspekte der Anwendung weiter zu benutzen. Aus den Konzepten und der Umsetzung ergeben sich viele weitere Ideen und Ansätze, die realisierbar sind und implementiert werden können.

Literaturverzeichnis

- [1] clips. CLIPS: A Tool for Building Expert Systems. <http://www.clipsrules.net/>.
- [2] David Embrey. The Skill, Rule and Knowledge Based Classification. *Human Error*, page 10.

Online-Quellen

- [3] Bedeutung AW (Arbeitswert). <https://wissen.autoxpert.de/hc/de/articles/360008973819-Bedeutung-AW-Arbeitswert->.
- [4] abbyy. PDF Software: Open, Read & Edit PDFs. <https://pdf.abbyy.com/>. [letzter Zugriff: 4. Juni 2022].
- [5] amazon. Intelligente Extraktion von Text und Daten mit OCR – Amazon Textract – Amazon Web Services. <https://aws.amazon.com/de/textract/>. [letzter Zugriff: 4. Juni 2022].
- [6] apache. Apache Tika – Apache Tika. <https://tika.apache.org/>. [letzter Zugriff: 4. Juni 2022].
- [7] atlanhq. Comparison with other PDF Table Extraction libraries and tools · atlanhq/camelot Wiki. <https://github.com/atlanhq/camelot>. [letzter Zugriff: 4. Juni 2022].
- [8] camelot. Camelot: PDF Table Extraction for Humans. Atlan Technologies Pvt Ltd, July 2022. [letzter Zugriff: 4. Juni 2022].
- [9] Adobe I/O-Adobe Developers. Extract Text from PDF | Extract Data from PDF | Visualizer - Adobe Developers. <https://developer.adobe.com/document-services/apis/pdf-extract/>. [letzter Zugriff: 4. Juni 2022].
- [10] Ben Dickson. What is symbolic artificial intelligence? <https://bdtechtalks.com/2019/11/18/what-is-symbolic-artificial-intelligence/>, November 2019. [letzter Zugriff: 4. Juni 2022].
- [11] Docsumo. PDF Scraper - Scrape data from pdf | PDF data extraction. <https://docsumo.com/blog/extract-data-from-pdf>, June 2022. [letzter Zugriff: 4. Juni 2022].

- [12] measures. Measuresforjustice/textricator: Textricator is a tool to extract text from documents and generate structured data. <https://github.com/measuresforjustice/textricator>. [letzter Zugriff: 4. Juni 2022].
- [13] pdfminer. Welcome to pdfminer.six's documentation! — pdfminer.six _VERSION_ documentation. <https://pdfminersix.readthedocs.io/en/latest/>.
- [14] Jeremy Singer-Vine. Pdfplumber, July 2022.
- [15] stichdata. What is Data Extraction? [Tools & Techniques]. <https://www.stitchdata.com/what-is-data-extraction/>. [letzter Zugriff: 4. Juni 2022].
- [16] tabulapdf. Tabulapdf/tabula: Tabula is a tool for liberating data tables trapped inside PDF files. <https://github.com/tabulapdf/tabula>. [letzter Zugriff: 4. Juni 2022].
- [17] Laura Theuerer. Was macht das After Sales Management im Autohaus? | Karriere-Ratgeber - kfz-betrieb Jobs. <https://jobs.kfz-betrieb.de/karriere-ratgeber/was-macht-das-after-sales-management-im-autohaus-51.html>. [letzter Zugriff: 4. Juni 2022].
- [18] unixuser. PDFMiner. <https://www.unixuser.org/~euske/python/pdfminer/>. [letzter Zugriff: 4. Juni 2022].