

06.19

KSI

Krisen-, Sanierungs- und Insolvenzberatung

Wirtschaft Recht Steuern

15. Jahrgang
November/Dezember 2019
Seiten 241–288

www.KSIdigital.de

Herausgeber:

Peter Depré, Rechtsanwalt und Wirtschaftsmediator (cvm), Fachanwalt für Insolvenzrecht

Dr. Lutz Mackebrandt, Unternehmensberater

Gerald Schwamberger, Wirtschaftsprüfer und Steuerberater, Göttingen

Herausgeberbeirat:

Prof. Dr. Markus W. Exler, Fachhochschule Kufstein

Prof. Dr. Paul J. Groß, Wirtschaftsprüfer, Steuerberater, Köln

WP/StB Prof. Dr. H.-Michael Korth, Präsident des StBV Niedersachsen/Sachsen-Anhalt e.V.

Dr. Harald Krehl, Senior Advisor, Wendelstein

Prof. Dr. Jens Leker, Westfälische Wilhelms-Universität Münster

Prof. Dr. Andreas Pinkwart, HHL Leipzig Graduate School of Management

Prof. Dr. Florian Stapper, Rechtsanwalt, Stapper/Jacobi/Schädlich Rechtsanwälte-Partnerschaft, Leipzig

Prof. Dr. Wilhelm Uhlenbruck, Richter a.D., Honorarprofessor an der Universität zu Köln

Prof. Dr. Henning Werner, Dekan der Fakultät für Wirtschaft, SRH Hochschule Heidelberg

Strategien Analysen Empfehlungen

Zielkaskadierung und Handlungsstrategien von Herstellern in der akuten Krise eines Automobilzulieferers [Dr. Alexander Jaroschinsky, 245]

Nach dem Crowdfunding-Erfolg in die Krise [Dr. Johannes Klein, 251]

Die Besteuerung von Sanierungsgewinnen [Prof. Dr. Sylvia Bös, 257]

Praxisforum Fallstudien Arbeitshilfen

Big Data und Künstliche Intelligenz: Praktischer Einsatz in Krisenunternehmen [Dr. Marcus Dill / Thomas Möllers, 263]

Der neue IDW S 2 zu Anforderungen an Insolvenzpläne [Jens Weber / Dr. Henrik Solmecke, 269]

Restrukturierungs- und Sanierungspraxis vor neuen Herausforderungen [Thesen-Kolloquium Restrukturierung/Sanierung vom 26. 9. 2019, 275]

Beilage

Jahresinhaltsverzeichnis 2019

Big Data und Künstliche Intelligenz: Praktischer Einsatz in Krisenunternehmen

Teil B: Case Study zum Einsatz von Big Data und Künstlicher Intelligenz

Dr. Marcus Dill und Thomas Möllers*

Mit der Suche, Identifikation, Analyse, Nutzung und Interpretation von relevanten Daten sind in Krisensituationen oft große Schwierigkeiten verbunden. Im Teil A dieses Beitrags in KSI 05/2019 S. 228–232 wurden die generelle Situation, die Grundlagen und die Voraussetzungen für den Einsatz von Big Data (BD)¹ und Künstlicher Intelligenz (KI) beschrieben. Im Teil B werden nun anhand einer Case Study mit einer realen Krisen- bzw. Restrukturierungs- bzw. Insolvenz-Situation innerhalb eines konkreten Projekts die Möglichkeiten und Herausforderungen von KI veranschaulicht². Dabei soll der Nachweis von Anfechtungs- und Schadensersatzforderungen mittels eines sog. eDiscovery³ detailliert beschrieben werden.

1. Einführung zu Teil B mit einem beispielhaften Praxisfall bei Anfechtungen und Schadensersatzforderungen

Die Case Study wird konkret belegen, dass im Rahmen eines ganzheitlichen Daten-Management Ansatzes KI-Technologien in Kombination mit klassischen Daten-Analysen in der Insolvenzverwaltung und in der Beratung die eigene und die Produktivität des Krisenunternehmens deutlich verbessern, indem

- Opportunities gesucht,
- Nutzen realisiert,
- Erlöse gesteigert,
- Kosten gesenkt und
- Bedrohungen vermieden

werden können. Beispielsweise wird die Case Study belegen, dass die Erweiterung des EDRM-Prozesses⁴ durch KI vielfältige positive Konsequenzen für die eDiscovery hatte. Die am Prozess beteiligten Personen verbrachten ihre Zeit nicht mit zähem Anlesen von Dokumenten, sondern wurden nach der Vorverarbeitung und Clusteranalyse systematisch an den Dokumentenstamm des Krisenunternehmens herangeführt. Dieses Big Picture aller Dokumente verbesserte und steuerte so das Data-Management. Auch die Ergebnisse der Begründungen von Anfechtungs- und Schadensersatzansprüchen gegenüber Dritten mittels KI zeigten Wirkung. Überwiegend konnte seitens der Insolvenzverwaltung auf Klagen verzichtet bzw. relativ schnell einvernehmliche Vereinbarungen getroffen werden, die die Masse in nicht unerheblichem Umfang vergrößerten⁵.

2. Kennzeichnung des Praxisfalls

Das im Praxisfall betrachtete Unternehmen verfügte über kein gutes Daten- und Dokumentenmanagement – typisch für Krisensituationen. Die Dokumente waren eher schlecht organisiert. Es war nicht möglich, für einen ISR-Anwendungsfall einen Teil der Dokumente automatisch zu exportieren.

Mithilfe eines sog. Early Data Assessment und eines sog. Data Culling⁶ konnte auf der Basis von Stichproben schnell ermittelt werden, dass die Exploration der gesamten Daten erfolgversprechend sein könnte. Daher wurde vom Insolvenzverwalter in Abstimmung mit dem Gläubigerausschuss der Einsatz von KI bzw. eDiscovery entschieden.

Die Daten des Rechnungswesens und der Materialwirtschaft wurden zunächst – obwohl zu Beginn des Insolvenzverfahrens noch in einer qualitativ eher schlechten Verfassung – durch ein Business-Management und mithilfe eines KI-Datenqualitätstools in einen verwertbaren Zustand versetzt, so dass sowohl die Vollständigkeit, Richtigkeit und Aktualität der Finanz- und der Bestandsbuchhaltung als auch die Qualität der Daten in einem ausreichenden Maße erreicht werden konnte.

Die IT-Systeme des Krisenunternehmens mit den relevanten Daten waren bereits bei Insolvenzeröffnung gesichert und in einem Archivrechenzentrum erfolgreich wiederhergestellt worden, was sich für die weitere Bearbeitung noch als einen großen Vorteil herausstellen sollte.

Aus Daten- und Geheimnisschutzgründen sind sowohl die originären Daten im Rahmen der Verarbeitung anonymisiert als auch die in dieser Veröffentlichung präsentierten Ergebnisse und Begrifflichkeiten pseudonymisiert.

* Dr. Marcus Dill, Geschäftsführer der mayato GmbH, E-Mail: marcus.dill@mayato.com. Dipl.-Kfm., M.Sc. Thomas Möllers, LL.M., Geschäftsführer der INSO Projects GmbH, E-Mail: thomas.moellers@inso-projects.de

1 Große Datenmengen mit den sog. 11 Firician Eigenschaften werden in Englisch als Big Data (BD) bezeichnet.

2 Das Projekt wurde in Kooperation zwischen der INSO Projects GmbH – schwerpunktmäßig tätig auf dem Gebiet des Daten-Managements in ISR – und der mayato GmbH – Experten für Textanalysen und Künstliche Intelligenz – durchgeführt.

3 eDiscovery ist ein Verfahren, bei dem für einen definierten Sachverhalt relevante Daten (zumeist in elektronischen Unterlagen wie E-Mails, E-Akten, Chat-Protokolle, PDF- und Office-Dokumenten) gesucht, identifiziert, analysiert, aufbereitet, interpretiert und bereitgestellt bzw. übergeben werden. Mittels definierter Prozesse müssen die Vollständigkeit, Richtigkeit und Aktualität der Daten sichergestellt und gleichzeitig die Gefahr, Geschäftsgeheimnisse zu verlieren, minimiert werden.

4 EDRM steht für Electronic Discovery Reference Model.

5 Die nachfolgende Case Study basiert auf dem sog. EDRM-Framework, das für ISR-Zwecke angepasst und um ein Business-, Projekt- und Daten-Management erweitert wurde. Eine Grafik zur Veranschaulichung des komplexen EDRM-Frameworks kann unter thomas.moellers@inso-projects.de angefordert werden.

6 Data Culling ist der Prozess der Suche und Isolierung von Daten basierend auf spezifischen Kriterien, wie z. B. Schlüsselbegriffe oder Zeiträume. Die drei gängigsten Methoden sind DeNISTing, Dedupe, Search terms (Suchbegriffe).

3. Durchführung des Datenmanagements

3.1 Datenextraktionen und Datenexport zwecks Datensatz-Erstellung

Im ersten Schritt wurden sehr große Datenextraktionen aus dem SAP ERP-System – insbesondere Rechnungswesen-Daten – mithilfe eines speziell für ISR-Zwecke entwickelten, leistungsfähigen IT-Tools für Datenextraktion, -analyse und -qualitätssicherung automatisiert durchgeführt.

In einem zweiten Schritt konnte der Datenexport aus DMS-⁷ und E-Mail-Systemen mit Hilfe spezieller Skripte und Programme schnell realisiert werden. Dazu mussten Forensiker auch gelöschte E-Mails und Office-Dateien wiederherstellen sowie geänderte Dokumente identifizieren und die gemachten Änderungen nachvollziehen.

Danach gelang es, die ERP-, DMS- und E-Mail-Daten zu einem Datensatz zusammen zu stellen. Damit mussten nicht jede Transaktion und jedes Dokument einzeln auf Relevanz geprüft werden, sondern es konnten durch sog. Meta-Informationen und Verzeichnisstrukturen eine Vielzahl an Transaktionen und Dokumenten gleichzeitig selektiert werden. Dieser Datensatz enthielt viele unterschiedliche Daten-/Dokumentenformate und -arten. Im vorliegenden Praxisfall traten zudem viele unterschiedliche Sprachen, Zeichensätze und Formate auf.

3.2 Analyse des Datensatzes

Für die Analyse des Datensatzes wurden zwei verschiedene Methoden eingesetzt:

(1) Mit **Clustering** wurden alle Dokumente anhand von Ähnlichkeiten neu organisiert. Das Ergebnis dieses Verfahrens ist eine saubere Einteilung aller Dokumente. Ein Mensch würde in diesem Schritt eine Ordnerstruktur erstellen und die Dokumente darin ablegen. Der Computer macht dies mit Clustering automatisch – insbesondere unter der Betrachtung des gesamten Datensatzes, was in dem Umfang für einen Menschen nicht möglich wäre. Die saubere Gruppierung der Dokumente unterstützte die vollständige Abgrenzung von Werten.

(2) Mit **Topic Models** wurde eine inhaltliche Analyse durchgeführt. Damit können gezielt bestimmte Themen aus den Dokumenten sys-

tematisch gebildet und anschließend extrahiert werden. Sind bestimmte Themen von hohem Interesse, können darüber die relevanten Dokumente gefunden werden.

3.3 Preprocessing

Vor der Anwendung von Algorithmen muss zunächst der Datensatz maschinenlesbar gemacht werden. Dies geschieht u. a. mit Hilfe von OCR-Software⁸, welche die PDF-⁹ und Office-Dateien in ein reines TXT-Format¹⁰ umwandelt. Die Qualität dieser Umwandlung ist äußerst wichtig für die Qualität der späteren Ergebnisse. Insbesondere sind systematische Fehler wie z. B. eine falsch erkannte Sprache zu vermeiden. Vereinzelt Fehler wie einmalig falsch erkannte Wörter hingegen stellen aufgrund ihrer probabilistischen Natur kein Problem für NLP-Algorithmen¹¹ dar und verschwinden in der Masse der Daten.

Sind die Daten im TXT-Format, können sie im nächsten Schritt vorverarbeitet werden. Da viele Preprocessing-Maßnahmen darauf basieren, dass die Sprache des Dokuments bekannt ist, werden die Dokumente zunächst nach Sprachen getrennt. Auch dies erfolgt automatisch mit Hilfe von leistungsfähigen NLP-Tools, die die Sprache eines Textes automatisch feststellen können.

Clustering und Topic Modelling arbeiten auf den einzelnen Wörtern eines Textes, weshalb ein nächster Preprocessing-Schritt die sog. **Tokenization** der Dokumente ist, bei dem u. a. Satzzeichen wie Kommata und Punkte von den Worten getrennt werden. Als Resultat liegt also nun jeder Text in Form einer Liste von Wörtern und wahlweise Satzzeichen vor, mit der im Folgenden gearbeitet wird.

3.4 Weitere Vorverarbeitungsschritte

Weitere sinnvolle Schritte der Vorverarbeitung sind das Herausfiltern von sog. **Stopwords**. Es handelt sich hierbei um Funktionswörter, wie Artikel, Bindewörter oder Präpositionen („der“, „und“, „in“), die keine Bedeutungsträger sind. Ihnen gegenüber stehen Inhaltswörter (Nomen, Verben, Adjektive und Adverbien). Hinzu kommt das Versehen der Wörter mit „Part-of-speech“-Tags, welche die Wortart anzeigen. Dieses sog. **POS-Tagging** ist selbst nicht trivial, da dasselbe Wort in Abhängigkeit von seinem

Kontext im Satz verschiedenen Wortarten angehören kann. Allerdings handelt es sich dabei um ein gut erforschtes Problem und inzwischen sind Systeme, die Wörter schnell und zuverlässig mit ihren POS-Tags versehen, einfach zu integrieren und für viele Sprachen verfügbar.

Auch das **Lemmatisieren** der Wörter ist in vielen Fällen hilfreich. Dabei wird ein Wort auf seine Grundform zurückgeführt, also beispielsweise ein Nomen im Plural auf den entsprechenden Singular oder ein flektiertes Verb auf seine Infinitivform. Waren also z. B. die Wörter „Person“ und „Personen“ für das System bislang völlig verschieden, werden sie nun beide als dasselbe Wort betrachtet. Dies hat Vorteile für die nachfolgenden Algorithmen.

Im letzten Vorverarbeitungsschritt müssen die Dokumente vergleichbar und deshalb auf dieselbe Art repräsentiert werden. Dafür wird für jedes Dokument ein sog. **tfidf-Vektor** erstellt. Dieser setzt sich zusammen aus der „term frequency“ („tf“) und der „inverse document frequency“ („idf“). Die **term frequency** gibt an, wie häufig ein Wort im betrachteten Dokument vorkommt. Da Wörter, die in sehr vielen Dokumenten vorkommen, wahrscheinlich weniger relevant für jedes einzelne Dokument sind (hiervon sind in erster Linie Funktionswörter betroffen), wird die term frequency anschließend mit der **inverse document frequency** multipliziert, welche umso niedriger ist, in desto mehr verschiedenen Dokumenten das Wort vorkommt.

Berechnet man die tfidf-Werte für jedes in der Dokumentensammlung vorkommende Wort und jedes Dokument, erhält man pro Dokument einen Vektor positiver rationaler Zahlen. Auch wenn dieser Vorverarbeitungsprozess zunächst recht aufwendig erscheint, wird er heutzutage vollständig und automatisiert von Computern übernommen. Das spezielle Know-how vorausgesetzt, reichen ein paar Zeilen Code, der Zugriff auf den Datensatz und eine entsprechende Rechenkapazität aus, um diese Vorverarbeitung rasch durchzuführen.

7 DMS steht für Document Management System.

8 OCR steht für Optical Character Recognition.

9 PDF steht für Portable Document Format.

10 TXT-Format ist ein reines Textformat.

11 NLP steht für Natural Programming Language.

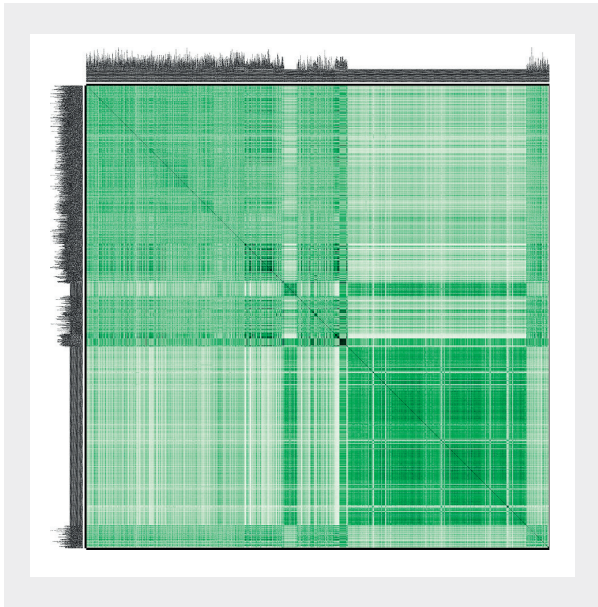


Abb. 1: Ähnlichkeitsmatrix aller Dokumente

Mit dieser Vorverarbeitung lassen sich die Ähnlichkeiten, die im Datensatz auftreten, bereits mit einfachen Verfahren darstellen.

Jede Zeile und Spalte steht in der Darstellung in Abb. 1 für ein Dokument. Insgesamt sind in diesem Beispiel etwa 3.600 Dokumente visualisiert, d. h. es werden 3.600 Zeilen und nochmals 3.600 Spalten belegt. Die farblichen Werte stehen für die Ähnlichkeit. Je grüner, desto ähnlicher sind die Dokumente. Zunächst fällt die Diagonale auf. Diese entsteht, weil dort das Dokument mit sich selbst verglichen wird. Der tf-idf-Vektor ist in dem Fall natürlich identisch, die Dokumente sind maximal ähnlich. Grundsätzlich sind in der Matrix unterschiedliche Ähnlichkeiten sichtbar. Diese Unterschiede sind die notwendige Voraussetzung für systematische Analysen.

Neben dieser Encodierungsart können die Dokumente heutzutage auch mittels Deep Learning unterteilt werden. Hierzu ist eine Reihe neuer Methoden veröffentlicht worden, die Computern ermöglichen, Dokumente nicht mehr in reinen Statistiken über die Wortanzahl darzustellen, sondern tatsächliches Textverständnis, Satzbau, Paragraphen und satzübergreifende Referenzen für die Encodierung zu berücksichtigen.

Zu diesem Durchbruch verhalf ein neuer Baustein in neuronalen Netzen: die sog. selbstre-

gulierte Aufmerksamkeit (engl. Self-Attention). Diese Komponente ermöglicht es neuronalen Netzen, explizit Fokus auf Zusammenhänge zwischen Wörtern zu legen, um auf sprachliche Feinheiten wie die Negation, Relativsätze oder den Bezug eines Personalpronomens zu erschließen.

Ein bekannter Vertreter an dieser Stelle ist BERT (Bidirectional Encoder Representations from Transformers), eine bei Google entwickelte Transformer-Architektur. Darin werden zunächst einzelne Wörter durch sog. Word Embeddings¹² repräsentiert. Ein neuronales Netz kriecht aus dieser Sequenz an Wörtern unter Berücksichtigung des

Kontextes (bidirektional) eine neue abstrahierte Repräsentation des Dokuments. Die Deep Learning Architektur ist mittlerweile auch fähig, deutsche Texte zu lesen und zu verstehen.

3.5 Clustering

Die aufbereiteten und vektorisierten Textdaten können nun für das Dokumenten-Clustering verwendet werden. Clustering eignet sich besonders gut als Lösungsansatz für sog. "unsupervised-learning"-Probleme, bei denen kaum Kennzahlen über die Verteilung der Daten existieren und explorativ

eine Ordnung gefunden werden soll. Ähnlichkeit zwischen Dokumenten kann nach verschiedenen Gesichtspunkten – wie Distanz-Metriken, Annahme einer Verteilung der Datenpunkte oder Dichte der Datenpunkte – definiert werden.

Grundsätzlich richten sich die Auswahl der Clusteralgorithmen und die Definition der benötigten Parameter nach der Größe und Beschaffenheit der Daten. Man unterscheidet beim Dokument-Clustering das sog. hierarchische Clustering und die Partitionierung. Der erste Ansatz ist besonders für kleinere Datenmengen geeignet, da jedes Dokument mit jedem anderen Dokument auf Ähnlichkeit geprüft wird, die Berechnungszeit also exponentiell mit der Anzahl der Dokumente steigt. Aus der Ähnlichkeitsberechnung entsteht ein Dendrogramm (Abb. 2), welches einer Baumstruktur ähnelt.

Für größere Dokumentenmengen ist der sog. K-Means-Algorithmus häufig eine gute Wahl. Er zählt zu der Gruppe der partitionierenden Clusteralgorithmen, wobei festgelegt werden muss, wie viele Cluster aus den Dokumenten geformt werden sollen und danach iterativ die optimalen Zentroiden dieser Cluster berechnet werden. In der folgenden Visualisierung repräsentiert jeder Punkt ein Dokument. Er wird hierfür anhand einer

¹² Word Embeddings sind zahlenbasierte Darstellungen einzelner Vokabeln in Form von Vektoren. Diese werden i. d. R. auf großen Wortschätzen und mit neuronalen Netzen gelernt. Hierbei ähneln sich Repräsentationen von Wörtern, die oft im selben Zusammenhang auftreten oder inhaltlich ähnlich sind.

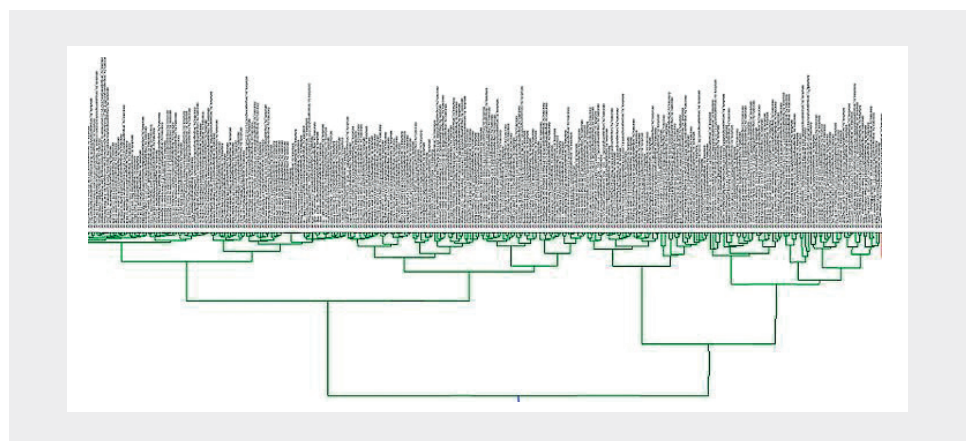


Abb. 2: Dendrogramm



Abb. 3: Visualisierung der Clusterzuteilung bei 30 Clustern (links) sowie 50 Clustern (rechts)

Dimensionsreduktion der tfidf-Vektoren ermittelt. Die Farbe steht für die Clusterzugehörigkeit des Dokuments. Indem jedes Dokument durch einen Punkt repräsentiert wird und Clusterzugehörigkeiten farblich ausgedrückt werden, sieht man, dass je nach Wahl der Clusteranzahl feinere Unterteilungen der Dokumente entstehen (vgl. Abb. 3).

Konkret konnte im Praxisfall anhand dieser Visualisierung eine Region interessanter Dokumente identifiziert werden. Die Verfeinerung der Clusterpartitionen half dabei, einen Überblick der Feinstruktur der Dokumente innerhalb dieser interessanten Region zu erhalten. An dieser Stelle wurden einige Zyklen des CRISP-DM-Prozesses durchlaufen. Sie verbesserten das Business Understanding (z. B. Relevanz für Anfechtungsansprüche), beeinflussten die Auswahl der Daten (Data Preparation) und verbesserten so das Clustering-Verfahren (Modelling).

3.6 Topic Models

Während die Dokumente beim Clustering vor allem nach Dokumenttypen gruppiert werden – also z. B. E-Mails, Verträge, Rechnungen, etc. – geht es beim Topic Modelling darum, die in den Dokumenten vorhandenen Themen zu extrahieren. Ein einzelnes Dokument kann dabei auch mehrere Themen enthalten und andererseits kann ein Thema in verschiedenen Dokumenttypen vorkommen.

Das hier verwendete Topic Modelling, also die Modellierung der vorhandenen Themen für eine Menge von Dokumenten, nennt sich „Latent Dirichlet Allocation“ (LDA) und fußt auf der Intuition, dass man sich das Entste-

hen eines Dokuments als einen generativen Prozess vorstellen kann, der auf den enthaltenen Topics und den dazugehörigen Wörtern basiert. Ein Topic Model besteht dementsprechend aus zwei Wahrscheinlichkeitsverteilungen.

Die erste Wahrscheinlichkeitsverteilung gibt an, aus welchen Wörtern sich ein Topic zusammensetzt. Die linke Tabelle der Abb. 4 zeigt beispielsweise ein Topic aus den Daten mit seinen zugehörigen Wörtern und ihren entsprechenden Wahrscheinlichkeiten. Als Betrachter kann man leicht erkennen, dass es sich bei diesem Topic um das Thema „Versicherungen“ handelt.

Die zweite Wahrscheinlichkeitsverteilung beschreibt die Verteilung der Topics innerhalb eines Dokuments. Enthält das Dokument nur ein Topic, hat dieses eine sehr hohe Wahrscheinlichkeit, während der Wert für die restlichen Topics annähernd null ist. Bei mehreren Themen im Dokument verteilt sich

die Wahrscheinlichkeit auf alle diese Themen. Die rechte Tabelle zeigt in Abb. 4 eine solche Verteilung von Topics auf Dokumente.

Die Gesamtanzahl aller Topics muss dabei im Vorhinein festgelegt werden. Außerdem verfügt LDA über einen Parameter, mit dem reguliert werden kann, wie hoch die Topic-Dichte innerhalb der einzelnen Dokumente sein soll. So können die Modelle noch genauer an die Eigenheiten der Daten angepasst werden.

Zur Bestimmung dieser beiden Wahrscheinlichkeitsverteilungen wird zunächst eine Dokument-Term-Matrix erzeugt, aus welcher dann unter Verwendung von Matrix-Faktorisierung die Topic-Wort-Matrix und die Dokument-Topic-Matrix entstehen.

Mit Hilfe von Sampling-Algorithmen wie dem sog. Gibbs Sampling werden diese beiden Matrizen iterativ verbessert, bis sie schließlich die finalen Wahrscheinlichkeitsverteilungen enthalten.

3.7 Auswertung

Möchte man nun also mit den Ergebnissen arbeiten, kann man zunächst in der Topic-Tabelle nach interessanten Themen suchen, um anschließend durch Sortieren und Filtern der Dokument-Topic-Matrix die entsprechenden Dokumente zu finden.

Die in Abb. 5 folgenden Wordclouds visualisieren einige der in den Daten gefundenen Themen anhand der zugehörigen Wörter, wobei die Größe eines Wortes von seiner Wahrscheinlichkeit im Topic abhängt.

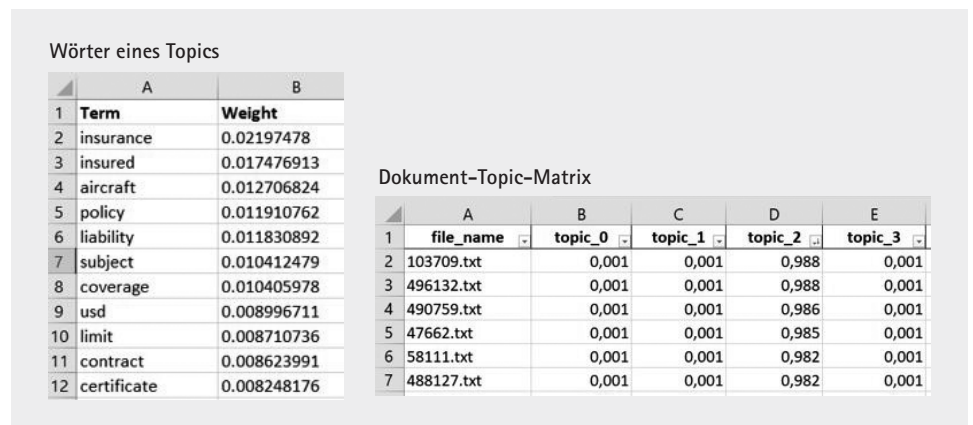


Abb. 4: Beispielhafte Ergebnisse für Themen und deren Vorkommen in den Dokumenten



Abb. 5: Wordclouds zu vier Themen

4. Ergebnisse im Praxisfall

4.1 Ergebnisse zu Anfechtungsansprüchen

Es konnte aus der riesigen Datenmenge eine Reihe von signifikanten Anfechtungstatbeständen exakt identifiziert werden. Erste Indikationen wurden zunächst mit Hilfe des Clustering ermittelt. So konnte eine Reihe von ähnlichen Vorgängen und Mustern in den Daten bzw. Dokumenten entdeckt werden. Mithilfe des Topic Models wurden die Inhalte analysiert und mittels Wordclouds ISR-relevante Themen und Aspekte in den Daten bzw. Dokumenten visualisiert. Dies ermöglichte unabhängig vom Dokumententyp anhand von Inhalten das Identifizieren von Anfechtungstatbeständen. Insbesondere konnte darüber eine

begrenzte Anzahl von konkreten Vorgängen in den E-Mail-Daten entdeckt werden, die über die dort enthaltenen Informationen mit Buchhaltungsdaten und Transaktionen verknüpfbar waren und damit belegbar wurden.

Die diesen Buchhaltungsdaten und Transaktionen zugehörigen gescannten Dokumente im DMS-System (z. B. Bestellungen, Lieferungen, Wareneingänge, Eingangsrechnungen etc.) und Zahlungs- und Verrechnungsvorgänge ließen sich über die Belegvernetzungen und die Auszifferungen im ERP-System ermitteln und fungierten damit als Nachweis. Es wurden so lückenlose Beweisketten geschaffen. Diese eigentlichen Nadeln im Heuhaufen ließen sich somit gut und schnell auffinden.

4.2 Ergebnisse zu Haftungsansprüchen aufgrund fehlerhaften Testats

Die materiellen Haftungstatbestände betreffen die Abschlussprüfung und die Testat-Gewährung durch Wirtschaftsprüfer. Die Indikation führte aber – im Gegensatz zu den Anfechtungsansprüchen – hier zunächst über die dynamische Analyse der Lagerbestände und der Bestandsbewertung. Mithilfe eines ISR-spezifischen Bestandsanalyseprogramms, das verschiedene statistische Verfahren umfasste, wurden zunächst die relevanten wesentlichen Objekte (Artikel, Geschäftspartner, Bestandstransaktionen) und Zeiträume bestimmt.

Ein dynamisches Clustering der betreffenden Rechnungswesens- und Materialwirtschaftsdaten ermöglichte dann die Identifikation von kritischen Mustern und Referenzen zu den relevanten Objekten im Zeitablauf. Es konnten damit die betreffenden Artikel, Buchungen und die zugehörigen gescannten Belege exakt ermittelt werden.

Es wurden anschließend – basierend auf den vorherigen Ergebnissen – korrespondierende Kommunikationsflüsse und -inhalte mit deutlichen Mustern innerhalb des Krisenunternehmens und mit den Abschlussprüfern in E-Mails durchsucht, was zu eindeutigen Hinweisen auf eine viel zu hohe Bestandsbewertung führte und eine darauf basierende Erteilung des Testats ermöglichte.

Zusammen mit den durch das Business Management gewonnenen Erkenntnissen konnte schnell nachgewiesen werden, dass bei genauer und sorgfältiger Prüfung des Rechnungswesens und der Materialwirtschaft das Testat in dieser Form hätte verweigert werden müssen.

U. a. kam es damit im Unternehmen zu einer formalen Verletzung der Ordnungsmäßigkeit bzw. der GoB. Auch hätten die Ergebnisse der Gewinn- und Verlustrechnung signifikant schlechter dargestellt werden müssen.

4.3 Effizienzvorteile

Die Erweiterung des EDRM-Prozesses durch KI hatte vielfältige positive Konsequenzen für die eDiscovery. Die am Prozess beteiligten Personen verbrachten ihre Zeit nicht mit zähem Anlesen von Dokumenten, sondern wurden nach der Vorverarbeitung und Clus-

teranalyse systematisch an den Dokumentenstamm des Krisenunternehmens herangeführt. Dieses Big Picture aller Dokumente verbesserte und steuerte so das Data-Management. Daraus folgten gezielte Anforderungen zu bestimmten ISR-Fragestellungen, die effizient und automatisch abgearbeitet werden konnten. Topic Models ermöglichten dabei die automatische Verknüpfung inhaltlich ähnlicher Dokumente für Digital Assets, Vertragsanalysen, Unternehmensplanung und rechtliche Prozesse. Die KI ist zudem kostengünstig und daher besonders für große Datenmengen geeignet. Eine händische, manuelle Verarbeitung wäre nicht wirtschaftlich gewesen.

Auch die Ergebnisse der Begründungen von Anfechtungs- und Schadensersatzansprüchen gegenüber Dritten mittels KI zeigten Wirkung. Überwiegend konnte seitens der Insolvenzverwaltung auf Klagen verzichtet bzw. relativ schnell einvernehmliche Vereinbarungen getroffen werden, die die Masse in nicht unerheblichem Umfang vergrößerten.

5. Ausblick auf weiterführende Analysen

Unter dem Aufbau eines Datensatzes für Gut- und Schlechtfälle wiederholt auftretender Muster in Krisenunternehmen können zudem neue KI-Systeme trainiert werden. Sie verstehen die Systematik aus vergangenen Daten und können so für neue Projekte automatisch und direkt die relevanten Dokumente identifizieren.

Basierend aus den Erfahrungen bereits vorhandener Projekte ermöglicht dieses Vorgehen auch eine Evaluierung des KI-Modells, z. B. über sog. Key Performance Indicators (KPI).

Der nächste Schritt bestünde nun in der Verwendung von weiterführenden KI-Techniken, die auf der Basis von bestehenden Klassifizierungen von Gut- und Schlecht-Fällen angelernt werden und im Anschluss Systematiken in Daten neuer Projekte automatisiert erkennen können. Dies erlaubt es dann, relevante Dokumente und Buchungen KI-gestützt sehr kurzfristig zu identifizieren. Voraussetzungen für dieses Vorgehen sind das Vorhanden-Sein einer entsprechenden Menge von klassifizierten Buchhaltungs-

Daten, Belegen und Informationen zu Prozessausgängen sowie die Zugriffsberechtigungen hierauf, wie sie fast ausschließlich bei IT-Dienstleistern für Insolvenzverwalter zu finden sein dürften.

6. Fazit

KI wird bereits in naher Zukunft die rechtliche, technische und betriebswirtschaftliche Krisenberatung auf bahnbrechende Weise transformieren. Sie bringt einen grundlegenden Wandel für die Durchführung aller Arten von Krisenprojekten und tangiert damit Insolvenzverwalter und Sanierungs- und Restrukturierungsberater gleichermaßen.

Kanzleien und Beratungen sind traditionell immer noch sehr papierlastig. Viele Arbeitsprozesse sind hier noch analog und mit einem sehr hohen Personal- und Papiereinsatz verbunden. Durch die smarte Verknüpfung von strukturierten ERP-Daten mit un- bzw. semi-strukturierten Daten aus Dokumenten-Management-Systemen (DMS) sowie aus E-Mails, Chats und Social Media eröffnen sich im Rahmen von KI völlig neue digitale Möglichkeiten für weitergehende Analysen, Entscheidungen und Aktionen.

Die Case Study hat eindrucksvoll aufgezeigt, welche enormen Potenziale und Möglichkeiten sich durch den Einsatz von KI ergeben können. Auch in Punkto Geschwindigkeit lassen sich damit deutliche Verbesserungen

für die Durchsetzung von Anfechtungs- und Schadensersatzforderungen erzielen. Konsequenter weitergedacht, sollte sich jeder Insolvenzverwalter bzw. Sanierungsberater die Fragen stellen, ob er zur Steigerung des Unternehmenserfolgs wirklich auf den Einsatz von KI verzichten will. Wie kann er wissen, ob die richtigen Beweismittel extrahiert worden und diese nun vollständig sind? Und kann hier ein Versäumnis vielleicht sogar eine Haftung nach sich ziehen?

Die Zukunft der Nutzung von KI in ISR wird auf jeden Fall spannend werden. Insgesamt gesehen werden mit dem KI-Einsatz umso mehr Freiräume und eine höhere Effektivität für noch anspruchsvollere Aufgaben wie die Sanierung von bzw. die Investorensuche für Krisenunternehmen geschaffen. KI ist damit ein entscheidender Wettbewerbs- und Überlebensfaktor in Krisensituationen geworden¹³.

¹³ Dem vorliegenden Beitrag liegt als Hauptquelle das Werk von Thomas Möllers über „Daten-Management in Krisenunternehmen“ aus dem Jahr 2017 zugrunde. Für den interessierten Leser sei zur Vertiefung ferner auf folgende Quellen verwiesen: Buxmann, Schmidt (Hrsg.), Künstliche Intelligenz, 1. Aufl. 2019; Coleman, Navigating the Labyrinth, 1. Aufl. 2018; DAMA International, Data Management Body of Knowledge (DMBOK), 2. Aufl. 2017; Freiknecht/Papp, Big Data in der Praxis, 2. Aufl. 2018; Gansor/Totok, Von der Strategie zum Business Intelligence Competency Center (BICC), 2. Aufl. 2015; Haneke et. al., Data Science, 1. Aufl. 2019; Hanschke et. al., Business Analyse, 2. Aufl. 2016; Michaeli, Competitive Intelligence, 1. Aufl. 2006.