

Chiara Zengerer

Nachvollziehbarkeit von Entscheiden – eine Gegenüberstellung von menschlichen Richterinnen und Richtern und künstlicher Intelligenz

Die Nachvollziehbarkeit von Ergebnissen stellt eine der grössten Herausforderungen beim Einsatz von künstlicher Intelligenz (KI) in der Justiz dar. Der vorliegende Beitrag untersucht, wie Entscheide einer KI im Vergleich zu Entscheiden von menschlichen Richterinnen und Richtern hinsichtlich ihrer Nachvollziehbarkeit zu beurteilen sind. Die Analyse diverser Beiträge zeigt, dass menschliche Entscheide nicht per se eine bessere Nachvollziehbarkeit bieten und dass den Gefahren einer KI teilweise mit Schutzmassnahmen und Kontrollen entgegengewirkt werden kann. Die grössten Hürden für eine KI bilden demgegenüber das fehlende Kontextwissen, die Unfähigkeit zur kritischen Reflexion sowie die Angabe der richtigen Gründe für einen Entscheid.

Beitragsart: Beiträge
Rechtsgebiete: LegalTech

Zitiervorschlag: Chiara Zengerer, Nachvollziehbarkeit von Entscheiden – eine Gegenüberstellung von menschlichen Richterinnen und Richtern und künstlicher Intelligenz, in: Jusletter IT 20. Juli 2023

Inhaltsübersicht

1. Einleitung
2. Automatisierung von rechtlichen Entscheiden
3. Problembereiche bei Entscheiden von künstlicher Intelligenz
 - 3.1. Wertungen und Abwägungen
 - 3.2. Unzureichende Begründungen
 - 3.3. Gefahr von Diskriminierungen
4. Vergleich zu menschlichen Entscheiden
 - 4.1. Einstellungen und Befindlichkeiten
 - 4.2. Unabhängigkeit und Unparteilichkeit
 - 4.3. Kontextwissen und Reflexion
5. Schlussfolgerungen

1. Einleitung

[1] Gemäss den Leitlinien «Künstliche Intelligenz» für den Bund stellt die Nachvollziehbarkeit von Ergebnissen eine der grössten Herausforderungen beim Einsatz von künstlicher Intelligenz (KI) dar. Das eidgenössische Departement für Wirtschaft, Bildung und Forschung legt fest, dass die Funktionsweise, der Zweck sowie die verwendeten Datensätze zum Training oder zur Entwicklung der KI in verantwortungsvoller und rechtskonformer Weise offengelegt werden sollen, damit die Entscheidungsprozesse für direkt und indirekt Betroffene nachvollziehbar und die Wirkungsweise für Fachleute überprüfbar sind.¹ Die inhaltliche Nachvollziehbarkeit hat in der Justiz eine besonders wichtige Bedeutung, da hoheitliche Entscheide den rechtsstaatlichen Grundsätzen genügen müssen.²

[2] Angesichts dieser fundamentalen Bedeutung und den dazu häufig vorgebrachten Einwänden gegen den Einsatz von KI drängt sich die Frage auf, wie berechtigt diese Einwände im Vergleich zu menschlichen Entscheiden sind. Zu diesem Zweck wird die Frage nach der technischen Machbarkeit von automatisierten Entscheiden sowie die Frage nach der positivrechtlichen Zulässigkeit beiseitegelassen. Der vorliegende Beitrag befasst sich allein mit der Nachvollziehbarkeit und stellt dabei konkret die Frage: Wie sind Entscheide einer KI im Vergleich zu Entscheiden von menschlichen Richterinnen und Richtern hinsichtlich ihrer Nachvollziehbarkeit zu beurteilen?

[3] Zur Einführung in die Thematik wird der Einsatz von KI zur Automatisierung von rechtlichen Entscheiden erläutert. Anschliessend werden im dritten Kapitel die Problembereiche bei Entscheiden von KI aufgezeigt, bevor im vierten Kapitel die untersuchten Problembereiche mit menschlichen Entscheiden verglichen werden. Abschliessend werden die Ergebnisse zusammengefasst und diskutiert.

2. Automatisierung von rechtlichen Entscheiden

[4] Der vorliegende Beitrag legt den Fokus auf den Einsatz von KI zur Automatisierung von rechtlichen Entscheiden. Eine Automatisierung der Entscheidungsfindung bedeutet, dass die Entscheide

¹ Eidgenössisches Departement für Wirtschaft, Bildung und Forschung WBF, Leitlinien «Künstliche Intelligenz» für den Bund, 2020, https://www.sbfi.admin.ch/dam/sbfi/de/dokumente/2020/11/leitlinie_ki.pdf.download.pdf.

² NINK DAVID, Justiz und Algorithmen, Über die Schwächen menschlicher Entscheidungsfindung und die Möglichkeiten neuer Technologien in der Rechtsprechung, Berlin 2021, S. 334.

autonom und somit ausschliesslich durch ein auf KI basierendes System getroffen werden.³ Hierfür besonders geeignet sind Verfahren, welche die Rechtsfolge an eindeutig messbare Kriterien binden.⁴ Bei Entscheiden, die im Einzelfall abgewogen werden müssen, besteht die Herausforderung vor allem in der Erkennung vorhandener Spielräume und ihrer ordnungsgemässen Nutzung.⁵

[5] Für die Automatisierung von rechtlichen Entscheiden sind vor allem Systeme maschinellen Lernens von Bedeutung. Die Funktionsweise solcher Systeme besteht darin, durch das Trainieren mit grösseren Mengen bestehender gerichtlicher Entscheide die relevanten Entscheidungsfaktoren zu erkennen und zu erlernen. Im Idealfall soll nach der Trainingsphase die Anwendung auch auf unbekannte Fälle möglich sein, indem das System auf das trainierte Entscheidungsmodell zurückgreift.⁶

[6] Die nachfolgend diskutierten Aspekte betreffen nur den Einsatz von KI in zivilrechtlichen Prozessen. In Strafprozessen wird die Datenerhebung durch Aussage- und Mitwirkungsverweigerungsrechte zusätzlich erschwert und die Festlegung eines Strafmasses beinhaltet weitere Faktoren wie Gerechtigkeitsempfinden und Resozialisierungschancen, die nur schwer algorithmisch erfasst werden können.⁷ Des Weiteren ist davon auszugehen, dass menschliche Richterinnen und Richter mindestens zur Feststellung des Sachverhalts unabdingbar bleiben, da technische Systeme bislang auf den Daten-Input durch Menschen angewiesen sind.⁸ Die Wahrheitsfindung durch eine KI ist ungleich schwieriger als die blossе Interpretation von Gesetzen auf der Basis von als wahr angenommenen Fakten.⁹

3. Problembereiche bei Entscheiden von künstlicher Intelligenz

[7] Das dritte Kapitel dieses Beitrags befasst sich mit den Einwänden hinsichtlich Nachvollziehbarkeit von Entscheiden, die durch eine KI gefällt wurden. Zu diesem Zweck werden die am häufigsten erwähnten Problembereiche beleuchtet.

3.1. Wertungen und Abwägungen

[8] Gerichtliche Entscheide beruhen regelmässig auf einer Abwägung der betroffenen Rechtsgüter auf dem Wege praktischer Konkordanz sowie auf Wertungs-, Beurteilungs- und Ermessensspielräumen. In diesen Fällen existieren verschiedene Begründungs- und Entscheidungsoptionen, die aus Sicht des Rechts alle «richtig» sind, aber trotzdem unterschiedlich gut vertretbar sein kön-

³ VON LUCKE JÖRN/ETSCHIED JAN, Wie Ansätze künstlicher Intelligenz die öffentliche Verwaltung und die Justiz verändern könnten, in: Jusletter IT 21. Dezember 2020, S. 258.

⁴ VON LUCKE/ETSCHIED, S. 263.

⁵ VON LUCKE/ETSCHIED, S. 259 f.

⁶ Zum Ganzen: DREYER STEPHAN/SCHMEES JOHANNES, Künstliche Intelligenz als Richter?, Wo keine Trainingsdaten, da kein Richter – Hindernisse, Risiken und Chancen der Automatisierung gerichtlicher Entscheidungen, Computer und Recht 11/35 2019, N1.

⁷ PUPPE FRANK, Gesellschaftliche Perspektiven einer fachspezifischen KI für automatisierte Entscheidungen, Informatik Spektrum 2/45 2022, S. 92.

⁸ NINK, S. 178.

⁹ PUPPE, S. 92.

nen.¹⁰ Bei solchen Entscheidungsspielräumen ist die Prüfung der Rationalität des Entscheidungsverfahrens sowie der Wissensgrundlagen für die Herleitung eines Entscheids besonders wichtig, was eine systemische Grenze für automatisierte Entscheide darstellen kann.¹¹

[9] GRECO sieht hingegen keinen Grund, der bereits im Vorhinein ausschliesst, einer künstlichen Intelligenz das juristische Werte beizubringen. Zu diesem Zweck könnte nach einer ersten Trainingsrunde mit bisherigen gerichtlichen Entscheiden in einer zweiten Kontrollrunde überprüft werden, ob die Ergebnisse der KI denen der Menschen entsprechen. Anschliessend würden die begangenen Fehler als Lernmaterial eingearbeitet werden. Nicht einmal Rechtsfortbildung liegt zwingend ausserhalb des Möglichen, da Rechtsfortbildung praktisch immer darin besteht, aus in anderen Rechtsgebieten bereits anerkannten Prämissen Folgen für einen anderen Sachverhalt abzuleiten und nicht darin, etwas wahrhaft Neues zu erschaffen.¹² Ist eine Rechtsfortbildung allerdings aufgrund eines Wertewandels in der Gesellschaft angezeigt, fehlt es am menschlichen Faktor, der die Anpassung der Rechtsprechung auslösen könnte.¹³

3.2. Unzureichende Begründungen

[10] Die Begründung ist im Sinne einer rationalitätsorientierten Rechtsauffassung für die Rechtmässigkeit des Entscheids konstitutiv. Die Begründung muss einerseits an sich nachvollziehbar sein und andererseits im Verhältnis zum Ergebnis nachvollziehbar sein. Bereits die erste Qualitätsanforderung erscheint nach dem heutigen Stand der Technik schwierig zu erreichen, da Algorithmen häufig Ergebnisse hervorbringen, ohne Auskunft darüber geben zu können, wie diese zustande gekommen sind. Einer KI die Fähigkeit beizubringen, nicht nur den richtigen Entscheid zu treffen, sondern auch die richtigen Gründe dafür anzugeben, erscheint jedoch nicht im Vorhinein als unmöglich, zumal juristische Texte häufig Textbausteine enthalten.¹⁴

[11] Für das Verhältnis der Begründung zum Ergebnis muss zunächst zwischen aufrichtigen Begründungen und Rationalisierungen unterschieden werden. Aufrichtige Begründungen stellen das Ideal dar, bei dem die Gründe für einen Entscheid auch die Motive für den Entscheid bilden. Als Rationalisierungen kritisiert werden Entscheide, die aus Motiven gefällt werden, die nicht zu Gründen werden können oder dürfen. Die Aufdeckung von Rationalisierungen ist allerdings bereits bei menschlichen Richterinnen und Richtern kaum möglich, da nicht in ihre Köpfe hineingeschaut werden kann. Die Unterscheidbarkeit ist sogar aussichtsreicher bei einer KI, da Programme entworfen werden könnten, die zusätzlich jeden Arbeitsschritt dokumentieren.¹⁵ Das Nachvollziehen einer solchen Dokumentation dürfte jedoch mit hohem Aufwand verbunden sein und das Ergebnis nur bei Routineentscheiden in typischen Fällen befriedigend ausfallen.¹⁶

¹⁰ Zum Ganzen: DREYER/SCHMEES, N 17.

¹¹ DREYER/SCHMEES, N 18.

¹² Zum Ganzen: GRECO LUIS, Richterliche Macht ohne richterliche Verantwortung – Warum es den Roboter-Richter nicht geben darf, RW Rechtswissenschaft 1/11 2020, S. 37 f.

¹³ WAGNER JENS, Legal Tech und Legal Robots, Wiesbaden 2020, S. 94.

¹⁴ Zum Ganzen: GRECO, S. 42 f.

¹⁵ Zum Ganzen: GRECO, S. 44 f.

¹⁶ GLESS SABINE/WOHLERS WOLFGANG, Subsumtionsautomat 2.0, Künstliche Intelligenz statt menschlicher Richter?, in: Böse Martin/Schumann Kay H./Toepel Friedrich (Hrsg.), Festschrift für Urs Kindhäuser zum 70. Geburtstag, Baden 2019, S. 159 f.

[12] Erstellt die KI ihren Entscheid auf der Basis einer Vielzahl historischer Entscheide mittels Mustererkennung, wird es im Normalfall keine eigenen, am Gesetz orientierten Begründungsschritte angeben können. Allerdings könnte die KI den historischen Parallelfall finden und die dortigen Begründungen mitliefern.¹⁷ Im Allgemeinen besteht jedoch die Gefahr von deutlich schematischeren und weniger ausdifferenzierten Entscheiden, da menschliche Richterinnen und Richter aufgrund der Lebenserfahrung Nuancen besser erkennen.¹⁸

3.3. Gefahr von Diskriminierungen

[13] Hinsichtlich der erhofften Objektivität von Entscheiden einer KI muss beachtet werden, dass ihre Entscheide nur so gut respektive so rational oder gerecht sein können wie die zugrundeliegenden Trainingsdaten.¹⁹ Das bedeutet, dass KI nicht von sich aus die notwendigen Schutzrechte für Schwächere oder Minderheiten gewährleistet, sondern dies Aufgabe des Staates bleibt.²⁰ Dieser Effekt wird zusätzlich verstärkt, wenn die gewonnenen Ergebnisse der KI ihrerseits als Grundlage für künftige Ergebnisse verwendet werden.²¹

[14] Das Argument der Diskriminierung kann jedoch lediglich eine Mahnung darstellen, die Möglichkeit der Diskriminierung ernst zu nehmen und ihr mit angemessenen technischen und organisatorischen Schutzmassnahmen und Kontrollen zu begegnen, da die Gefahr der Diskriminierung auch bei menschlichen Richterinnen und Richtern nicht ausgeschlossen werden kann. Es ist eher fraglich, ob sie ihnen nicht in einem noch stärkeren Masse ausgeliefert sind. Bei von einer KI gewonnenen Entscheiden kann zumindest mit Sicherheit ausgeschlossen werden, dass diese Ergebnisse planmässig herbeigeführt wurden.²²

4. Vergleich zu menschlichen Entscheiden

[15] In diesem Kapitel werden die oben beleuchteten Problembereiche bei Entscheiden einer KI mit den Problembereichen bei menschlichen Entscheiden verglichen. Dazu werden ebenfalls die am häufigsten vorgebrachten potenziellen Problemfelder untersucht und einer KI gegenübergestellt.

4.1. Einstellungen und Befindlichkeiten

[16] Wo Spielraum für Wertungen besteht, tragen menschliche Richterinnen und Richter immer auch persönliche Elemente in die Wahrnehmung des Sachverhalts hinein.²³ Unser Rechtssystem ist offen und bisweilen auch angewiesen auf den Einbezug des gesellschaftlichen Kontexts und

¹⁷ Zum Ganzen: WAGNER, S. 91.

¹⁸ WAGNER, S. 93 f.

¹⁹ NINK, S. 167.

²⁰ NINK, S. 168.

²¹ GRECO, S. 39.

²² Zum Ganzen: GRECO, S. 40 f.

²³ NINK, S. 41 f.

der aktuellen Wertmassstäbe in die Auslegung und Anwendung von Normen.²⁴ Ein Entscheid stützt sich deshalb typischerweise immer auch auf intuitiv-wertende Erkenntnisse. Allerdings dürfen einem Entscheid nicht diese «Gefühle», sondern nur die Fakten sowie ihre Gewichtung zugrunde gelegt werden, da ansonsten die Gefahr der Willkür durch übermässige Subjektivität und strukturell bedingte Ungleichbehandlung gleicher Sachverhalte entsteht.²⁵

[17] Im Gegenzug dazu würde eine KI unter gleichen Bedingungen vorhersagbar die gleichen Entscheide treffen, frei von Voreinstellungen und aktuellen Befindlichkeiten. Damit ginge eine grössere Transparenz und Standardisierung einher. Gleichzeitig besteht jedoch der Nachteil, dass von subtilen Besonderheiten jedes Einzelfalles abstrahiert werden muss, da abstrakte Begriffe zwar immer weiter konkretisiert werden können, aber nicht alle Konstellationen vorhersehbar sind. Daher müsste eine Berufung an eine höhere Instanz möglich sein, bei welcher die Entscheide von menschlichen Richterinnen und Richtern gefällt werden.²⁶

4.2. Unabhängigkeit und Unparteilichkeit

[18] Die Fähigkeit des Menschen, eine ganzheitliche Sichtweise einzunehmen, öffnet auf der Kehrseite ein mögliches Einfallstor für subjektive Einschätzungen, die sich vielleicht nur aus der individuellen Biografie oder einem informellen Entscheidungsprogramm erklären lassen.²⁷ In diesem Zusammenhang könnten Maschinen in viel grösserem Umfang unabhängig und unparteilich sein als Menschen, sofern das Programm immer wieder korrigierend kalibriert wird, um den bereits im Trainingsdatensatz enthaltenen Vorurteilen entgegen zu wirken.²⁸

[19] Bei menschlichen Richterinnen und Richtern kann nie ausgeschlossen werden, dass sie befangen sind, sich von aussen unter Druck setzen lassen oder sonst einen unsachlichen Entscheid treffen. Beim Einsatz von KI entstehen allerdings dadurch, dass sich schon ein einzelner unzulässiger Eingriff auf eine Vielzahl von Entscheiden auswirken kann, neue Angriffspotentiale. Ebenfalls könnten zentral entwickelte und eingesetzte Algorithmen von politischen Akteuren missbraucht werden.²⁹

4.3. Kontextwissen und Reflexion

[20] Von menschlichen Richterinnen und Richtern wird erwartet, dass sie bei einem Entscheid das Ergebnis in einer holistischen Gesamtbewertung hinterfragen und so reflektiert begründen, dass es intersubjektiv nachvollziehbar ist. Dazu ist die Fähigkeit zur offenen und kritischen Reflexion des eigenen Entscheids nötig. GLESS/WOHLERS sind der Meinung, dass nur Menschen in der Lage sind, ihren Entscheid adäquat zu begründen, Korrekturbedarf zu erkennen und so dem Berufsethos gerecht zu werden.³⁰

²⁴ NINK, S. 42.

²⁵ Zum Ganzen: NINK, S. 43 f.

²⁶ Zum Ganzen: PUPPE, S. 93.

²⁷ GLESS/WOHLERS, S. 160 f.

²⁸ GLESS/WOHLERS, S. 162 f.

²⁹ Zum Ganzen: WAGNER, S. 93.

³⁰ Zum Ganzen: GLESS/WOHLERS, S. 159.

[21] Für KI besteht auf der semantischen Ebene zusätzlich die Herausforderung, dass teilweise gleiche Sinngehalte mit unterschiedlichen Wörtern und teilweise unterschiedliche Sinngehalte mit ähnlichen oder gleichen Wörtern beschrieben werden können. Um die Bedeutung dieser unterschiedlichen Sinngehalte und damit auch den Kontext selbst bei simplen sozialen Komplexen zu verstehen, ist ein umfassendes Weltverständnis nötig. Verfahren maschinellen Lernens agieren jedoch bei ihren Wahrscheinlichkeitsrechnungen primär anhand identifizierter Korrelationen von mehreren Variablen. Sie verfügen damit nicht über das erforderliche Kontextwissen oder moralische Bewusstsein.³¹

5. Schlussfolgerungen

[22] Die Analyse diverser Beiträge zeigt, dass hinsichtlich der Nachvollziehbarkeit beim Einsatz von KI zur automatisierten Entscheidungsfindung im Vergleich zu menschlichen Entscheiden Uneinigkeit herrscht. Bei der Gegenüberstellung von KI und menschlichen Richterinnen und Richtern wird jedoch deutlich, dass menschliche Entscheide nicht per se eine bessere Nachvollziehbarkeit im Bereich der vorgebrachten Einwände gegen KI bieten.

[23] Der Gefahr der Diskriminierung durch eine KI kann mit angemessenen technischen und organisatorischen Schutzmassnahmen und Kontrollen sowie der sorgfältigen Auswahl der Daten begegnet werden. Zudem kann mit Sicherheit ausgeschlossen werden, dass solche Ergebnisse planmässig herbeigeführt wurden, was bei menschlichen Richterinnen und Richtern nur gehofft werden kann. Ebenso könnte eine KI in viel grösserem Masse als Menschen unabhängig und unparteilich sein.

[24] Auf der anderen Seite stehen das fehlende Kontextwissen einer KI und die Unfähigkeit zur kritischen Reflexion des eigenen Entscheids im Sinne einer holistischen Gesamtbewertung. Bei Wertungs-, Beurteilungs- und Ermessensspielräumen ist die Prüfung der Rationalität des Entscheidungsverfahrens und der Wissensgrundlagen für die Herleitung eines Entscheids besonders wichtig. Dies kann eine systemische Grenze für automatisierte Entscheide darstellen, solange einer KI das juristische Werten nicht beigebracht werden kann.

[25] Da die Begründung für die Rechtmässigkeit eines Entscheids konstitutiv ist, muss einer KI die Fähigkeit beigebracht werden, die richtigen Gründe für einen Entscheid anzugeben. Selbst wenn es gelingt, ein Programm zu entwerfen, das jeden Arbeitsschritt dokumentiert, dürfte das Nachvollziehen einer solchen Dokumentation mit hohem Aufwand verbunden sein. Im Allgemeinen besteht die Gefahr von deutlich schematischeren und weniger ausdifferenzierten Entscheiden.

[26] Die fortschreitende Digitalisierung in den Gerichten bietet die Grundlage für eine mögliche technische Umsetzung. Neben der vorliegend beiseitegelassenen Frage nach der technischen Machbarkeit von automatisierten Entscheiden sowie die Frage nach der positivrechtlichen Zulässigkeit ist ein öffentlicher Diskurs hinsichtlich der Wünschbarkeit solcher Systeme sowie die Präzisierung der Anforderungen hinsichtlich der Nachvollziehbarkeit in Bezug auf KI unabdingbar.

³¹ Zum Ganzen: DREYER/SCHMEES, N 15.

CHIARA ZENGERER ist Bachelor-Studentin der Rechtswissenschaften an der Universität St. Gallen (HSG).