

Janice Deborah Kowalski

Die Vorhersage von Gerichtsurteilen des EGMR durch Machine Learning – eine Fallstudie

Machine Learning Methoden werden heutzutage in immer mehr Bereichen eingesetzt. Auch in der Jurisprudenz gewinnen diese Methoden an Bedeutung. Die vorliegende Arbeit befasst sich mit der Vorhersage von Gerichtsurteilen des Europäischen Gerichtshofs für Menschenrechte anhand der Anwendung von Machine Learning, dabei wird die Methode des Natural Language Processing in Kombination mit Klassifikationsmethoden auf ihre Prädiktionsfähigkeit untersucht. Dieser Beitrag zeigt schlussendlich, wie Machine Learning im Schweizerischen Rechtssystem eingesetzt werden könnte, um Gerichtsent-scheide vorauszusagen.

Beitragsart: Next Generation

Region: Schweiz; EU

Rechtsgebiete: Artificial Intelligence & Recht; Rechtsinformatik; LegalTech

Zitiervorschlag: Janice Deborah Kowalski, Die Vorhersage von Gerichtsurteilen des EGMR durch Machine Learning – eine Fallstudie, in: Jusletter IT 30. September 2021

Inhaltsübersicht

1. Einleitung
2. Anwendungshintergrund
3. Experimente
 - 3.1. Erstes Experiment
 - 3.2. Zweites Experiment
 - 3.3. Drittes Experiment
4. Diskussion
5. Fazit

1. Einleitung

[1] Die Zunahme der Digitalisierung in der Jurisprudenz widerspiegelt sich sinnbildlich an der steigenden Anzahl von digital erhältlichen rechtlichen Dokumenten.¹ Dies gilt BATSAKIS et al. zur Folge insbesondere dort, wo es um juristische Argumentationen geht: «[...] *legal reasoning was one of the first application domains of artificial intelligence since its emergence, and more so after the proliferation of expert systems during the 1980s*».² Mit dieser Entwicklung, Gerichtsentscheide öffentlich zugänglich zu machen, wurde die Grundlage für «*Big Data*» Analysen bzw. Machine Learning Modelle geschaffen.³ Machine Learning gilt als Unterkategorie von LegalTech, die «*Verbindung von Recht und Technologie*»⁴, wobei Machine Learning wiederum der Technologie 2.0 angehört.⁵ Bei Letzterem werden technische Modelle oder Lösungen als Ersatz menschlicher Arbeit verwendet. So werden heutzutage Machine Learning Tools als Hilfestellung für oben genannte Problematik in der Jurisprudenz eingesetzt, wie bei der Vorhersage der Rechtsprechung. Dementsprechend analysiert die vorliegende Arbeit die Studie von MEDVEDEVA, VOLS und WIELING zur Vorhersage von Gerichtsurteilen des Europäischen Gerichtshofs für Menschenrechte (EGMR) anhand von Machine Learning. Es soll aufgezeigt werden, inwiefern sich das *Natural Language Processing* in Kombination mit *Klassifikationsmethoden* in den Rechtswissenschaften als zukunfts-fähig erweist.

2. Anwendungshintergrund

[2] Während sich in den Vereinigten Staaten solche quantitativen Analysen schon länger etabliert haben, sind anderenorts vorwiegend sogenannte *doctrinal research methods* anzutreffen. Obwohl verhältnismässig wenig rechtliche Studien auf der Anwendung von Machine Learning basieren, nimmt diese Entwicklung in den letzten Jahren international allmählich zu. So veröffentlichten ALETRAS et al. im Jahr 2016 eine wissenschaftliche Arbeit zur gleichen Thematik wie die der vorliegenden Arbeit.⁶ Das oben erwähnte Paper von MEDVEDEVA, VOLS und WIELING baut auf Letzterem auf, verwendet aber eine grössere Datenmenge bzgl. der Anzahl analysierten Gesetzesartikeln und den dazugehörigen Gerichtsentscheiden der Europäischen Menschenrechtskonvention

¹ GLASER/SCEPENKOVA/MATTHES, S. 61.

² BATSAKIS et al., S. 31.

³ MEDVEDEVA/VOLS/WIELING, S. 237.

⁴ KUMMER/PFÄFFLI, S. 134.

⁵ KUMMER/PFÄFFLI, S. 135.

⁶ MEDVEDEVA/VOLS/WIELING, S. 239 ff., vgl. ALETRAS et al.

(EMRK).⁷ Die vorliegende Studie basiert nur auf bereits entschiedenen und nicht künftigen Fällen des EGMR.

3. Experimente

[3] In der Studie von MEDVEDEVA, VOLS und WIELING wurden drei verschiedene Experimente mit der Anwendung von *Natural Language Processing* und der *Klassifikation* durchgeführt, welche nachfolgend analysiert werden. Bei allen Experimenten wurde das Modell jeweils trainiert und getestet, beim dritten Experiment geschah dies jedoch mit unvollständigen Informationen.⁸

3.1. Erstes Experiment

[4] Im ersten durchgeführten Experiment der Studie wurden Wörter und Sätze der verwendeten Fälle als Input genutzt, um die Gerichtsentscheide vorherzusagen. Als Grundlage dafür dienten alle englischen und auf der Webseite HUDOC des EGMR verfügbaren Entscheide bis zum Stand vom 11. September 2017. Das Ziel dessen stellte das Einteilen des jeweiligen Falls in zwei Kategorien dar, das heisst, es sollte jeder Fall entweder als *Verstoss* oder als *kein Verstoss* des vorliegenden Gesetzesartikels der EMRK eingestuft werden. Aufgrund dessen sollten die Entscheide der Richter dazu vorhergesagt werden.⁹ Als Methode wurde vom Machine Learning Bereich *Natural Language Processing* Gebrauch gemacht. Die *Klassifikationsmethode* kann grundsätzlich einen Text per se nicht einordnen, daher wurden die Texte in für die Maschine lesbare *n-grams* umgewandelt. Diese *uni-, bi- oder trigrams* wurden je nach Bedeutung für den Fall gewichtet und als Zahlen dargestellt. Die hier verwendete *Methode der Klassifikation* wird als *supervised* Machine Learning bezeichnet.¹⁰

[5] Die *Klassifikation* des Outputs in separate Kategorien ist als ein Teilbereich des Machine Learnings einzustufen.¹¹ In der vorliegenden Studie wurde dafür der Ansatz der *Support Vector Machine (SVM)* eingesetzt. Es handelt sich dabei um ein binäres Modell. Für jeden Gesetzesartikel wurde das Modell einzeln trainiert, damit dieses die zugehörigen Gerichtsfälle als *Verstoss* oder als *kein Verstoss* der Gesetzesartikel der EMRK klassifiziert.¹² Um das Modell für beide Klassen gleichmässig zu trainieren, wurden die Datensätze ausgewogen. So wurden von beiden Kategorien eine identische Anzahl Fälle verwendet.¹³ Allerdings liessen MEDVEDEVA, VOLS und WIELING die rechtliche Argumentation im Modell explizit aussen vor, sodass die Vorhersagen nur auf Fakten und nicht auf für den Verfahrensausgang richtungsweisenden Argumentationen und Gesetzesartikeln basierten. Die Aufteilung der Daten wurde dann zufällig in einen *training*

⁷ Vgl. MEDVEDEVA/VOLS/WIELING, S. 242: Analyse von neun Gesetzesartikeln der EMRK; ALETRAS et al., S. 6: Analyse von drei Gesetzesartikeln der EMRK.

⁸ Vgl. MEDVEDEVA/VOLS/WIELING, S. 248, 259.

⁹ MEDVEDEVA/VOLS/WIELING, S. 242 f., 246.

¹⁰ Vgl. SCHRIDER/KERN, S. 301 ff.: *Supervised* Machine Learning beruht auf bereits vorhandenem Wissen eines Musterdatensatzes, was für Vorhersagen für neue Daten verwendet wird im Unterschied zu *unsupervised* Machine Learning.

¹¹ HASTIE/TIBSHIRANI/FRIEDMAN, S. 101 ff.

¹² MEDVEDEVA/VOLS/WIELING, S. 243 f.

¹³ MEDVEDEVA/VOLS/WIELING, S. 247 f.

und *testing* Teil gegliedert. In der Phase des *trainings* liess man das Modell den Zusammenhang zwischen den Texten und den Gerichtsentscheiden lernen. Hingegen in der Phase des *testings* schätzte das Modell eigenständig die Entscheide anhand von den Texten alleine. Die Evaluierung wurde schliesslich anhand vom Anteil der korrekt vorhergesagten Entscheide vorgenommen. Hierfür wurden die Vorhersagen des Modells mit den tatsächlichen Entscheiden aus dem *testing* Teil abgeglichen.¹⁴ Um die Stabilität der Ergebnisse zu prüfen, wurde die Aufteilung in den *training* und *testing* Teil, als auch die Schätzung des Modells, mehrmals wiederholt. Dies nennt sich *cross-validation*.¹⁵

[6] Es bleibt bzgl. den Resultaten vom Experiment eins festzuhalten, dass durchschnittlich Vorhersagen von Gerichtsentscheiden des EGMR mit einer Präzision von ungefähr 75 % bei beiden Evaluationskriterien erreicht wurden. Würde man die Resultate zufällig schätzen, läge die Wahrscheinlichkeit bei 50 %. Dementsprechend kann auf die Prädiktionsfähigkeit der Methode geschlossen werden, obwohl der Input lediglich auf Textbausteinen der jeweiligen Gerichtsfälle beruht. Dies wurde auch mittels zweier weiterer statistischer Kriterien bestätigt.¹⁶ Jedoch können Unterschiede in der Genauigkeit der Vorhersage je nach Gesetzesartikel festgestellt werden. Dabei gilt generell, je mehr Daten zur Verfügung stehen, desto grösser ist die Präzision bzgl. ähnlicher Fälle zu einem Gesetzesartikel.

3.2. Zweites Experiment

[7] Beim zweiten Experiment wurde im Vergleich zu Experiment eins nichts abgeändert, ausser der Wahl der zeitlich chronologischen Reihenfolge der Daten für das *training* und *testing* des Modells, um effektiv zukünftige Gerichtsfälle vorherzusagen.¹⁷ Bei Experiment eins wurde eine solche Abfolge nicht eingehalten. Experiment zwei hat den Vorteil, dass diese Vorgehensweise logischer und realitätsnäher ist. Gleichwohl mussten die Gerichtsfälle neu balanciert werden, um die zeitliche Chronologie einhalten zu können. Die Konsequenz davon waren Datenverluste, wobei MEDVEDEVA, VOLS und WIELING sich aufgrund dessen nur noch auf drei der ursprünglich neun Gesetzesartikel der EMRK konzentrierten (Artikel 3, 6 und 8 der EMRK). Bei diesen verbleibenden Gesetzesartikeln blieben genügend Daten übrig, damit die Aussagekraft der Vorhersage nicht verloren geht.¹⁸ Schliesslich konnte hierbei festgestellt werden, dass das Vorhersagen von künftigen Gerichtsentscheiden des EGMR basierend auf vergangenen Fällen schwieriger war und die Präzision dessen sank. Grundsätzlich kann gefolgert werden, je grösser der zeitliche Abstand zwischen bereits ergangenen und zukünftigen Gerichtsurteilen, desto unpräziser ist die Vorhersage.

¹⁴ MEDVEDEVA/VOLS/WIELING, S. 243, 248 f.

¹⁵ MEDVEDEVA/VOLS/WIELING, S. 244.

¹⁶ Vgl. MEDVEDEVA/VOLS/WIELING, S. 253 f.

¹⁷ MEDVEDEVA/VOLS/WIELING, S. 257 f.

¹⁸ MEDVEDEVA/VOLS/WIELING, S. 257 f.

3.3. Drittes Experiment

[8] Experiment drei weicht von den ersten beiden Experimenten ab, so wurden die Nachnamen der jeweils beteiligten Richter als einziger Input verwendet. Die eigentlichen Textelemente der Urteile und somit auch deren Gewichtung waren kein Bestandteil dieser Vorhersage. Als ausschlaggebend galt lediglich, ob ein Richter einer Kammer in einem bestimmten Entscheid urteilte oder nicht.¹⁹

[9] Grundsätzlich wiesen die Resultate dieser Schätzung eine niedrigere Genauigkeit auf (durchschnittlich 65 %), wie wenn die Textteile von Entscheiden als Input genutzt wurden (durchschnittlich 75 %). Dies erscheint nachvollziehbar, zumal der Informationsgehalt bei einem Nachnamen kleiner war, wie bei einem Text. Allerdings zeigte die Präzision je nach Gesetzesartikel der EMRK eine grosse Heterogenität auf. So schätzte das Modell bei gewissen Artikeln ebenso genau wie bei einem Text. Bei anderen Gesetzesartikeln wurde nur leicht besser vorhergesagt, wie wenn zufällig geraten würde.²⁰

[10] Die Ergebnisse in Bezug auf die Bedeutung der massgebenden Wörter lassen sich in Experiment drei eindeutiger interpretieren als in Experiment eins. Es lässt sich also feststellen, welche Richter häufiger in einer Kammer dabei waren, welche für einen *Verstoss* oder für *keinen Verstoss* eines Gesetzesartikel der EMRK gestimmt hatten.²¹ Allerdings kann von keinem kausalen Zusammenhang ausgegangen werden. So besteht keine Annahme, dass Richter des EGMR immer pauschal für oder gegen einen Verstoss der EMRK urteilen würden.

4. Diskussion

[11] In den vorhergehenden Kapiteln konnte gezeigt werden, dass die vorliegende Studie von MEDVEDEVA, VOLS und WIELING die Vorstudie von ALETRAS et al. weiterentwickelte. MEDVEDEVA, VOLS und WIELING basierten ihre Ergebnisse zum einen auf einer grösseren Anzahl Gesetzesartikel sowie auf mehr Gerichtsurteilen pro Gesetzesartikel. Zum anderen liessen sie rechtliche Argumentationen bewusst aussen vor, um dem Modell nicht im Vorfeld Hinweise auf den endgültigen Entscheid zu liefern. Letzteres wurde von der Vorstudie nicht beachtet, was unter Umständen deren Ergebnisse verfälscht hatte. Der Chronologie der Reihenfolge in Bezug auf das *training* und *testing* hatte die Vorstudie keine Beachtung geschenkt.²² Dies erklärt auch, weshalb die Genauigkeit der Vorhersage bei ALETRAS et al. bei ungefähr 80 % und diese bei der vorliegenden Studie lediglich bei 75 % lag.²³

[12] Grundsätzlich können einige Vorteile der Studie von MEDVEDEVA, VOLS und WIELING für Vorhersagen im Allgemeinen genannt werden. So handelt es sich bei der hierfür verwendeten Methode der *Support Vector Machine* um keine sogenannte Black Box Methode, das heisst, der Algorithmus kann intuitiv interpretiert werden. Damit kann einerseits erklärt werden, welche Begriffe für die Vorhersagen von Bedeutung waren. Andererseits deutet dies auf keinen kausalen

¹⁹ MEDVEDEVA/VOLS/WIELING, S. 259 f.

²⁰ MEDVEDEVA/VOLS/WIELING, S. 259 f.

²¹ MEDVEDEVA/VOLS/WIELING, S. 261.

²² MEDVEDEVA/VOLS/WIELING, S. 259.

²³ Vgl. ALETRAS et al., S. 2; MEDVEDEVA/VOLS/WIELING, S. 237.

Korrelationszusammenhang hin. Ferner ist die Komplexität des *Natural Language Processing* verhältnismässig, sodass sie auf einfache Art und Weise repliziert und auf andere Fälle angewandt werden kann. Jedoch fehlt es teils an inhaltlicher Bedeutung der *n-grams*, welche die Textbausteine der Urteile repräsentierten.²⁴ Diese Problematik könnte z.B. mit semantischen Analysen verhindert werden.²⁵

[13] Gleichermassen sind gewisse Nachteile der Methode von MEDVEDEVA, VOLS und WIELING zu erwähnen. Aufgrund des Balancierens der beiden Kategorien *Verstoss* oder *kein Verstoss* der Gesetzesartikel der EMRK wurden Daten ausgelassen, was zu einer schlechteren Performance der Methode führte. Es lässt sich folgern, dass die Präzision der Methode trotz dem hohen Innovationsgehalt noch ausbaufähig ist. Mit *under-* oder *oversampling* hätte dem entgegengewirkt werden können, da mit diesen Methoden keine Datenmengen aufgegeben werden müssen.²⁶ Als weiterer Kritikpunkt gilt, dass technische Aspekte im Vergleich zu den rechtlichen in der Studie mehr Gewicht erhielten. So wurde nicht direkt beleuchtet, was solche Vorhersagen für die Jurisprudenz bedeuten würden, ob sie überhaupt rechtlich legitimiert sind und wo der EGMR oder natürliche Personen konkret von der Methode Gebrauch machen könnten. Dies könnte an der noch verbesserungsfähigen Genauigkeit der Vorhersagen oder an der Simplizität der Methode liegen. Mit anderen komplexeren Machine Learning Methoden, wie z.B. das Deep Learning, könnte dem entgegengewirkt werden.²⁷

[14] Es wurde von den Autoren jedoch nie beabsichtigt die Richter des EGMR zu ersetzen, sondern nur die Erfolgchancen eines Gerichtsfalles einschätzen zu können. Dies dürfte vor allem natürlichen Personen dienen. Dementsprechend wurden die Erkenntnisse der Studie mit der Webseite «JURI says» einer breiten Öffentlichkeit zugänglich gemacht, wo richtige und falsche Vorhersagen des Modells einsehbar sind.²⁸

[15] Gegebenenfalls könnte die Methode alternativ bei der Prüfung der Zuständigkeit eines Gerichts eingesetzt werden.²⁹ So könnte das Modell konkret die Arbeit in den Gerichten erleichtern, da die Entscheidung der Zuständigkeitsprüfung auch nicht so weitreichend wäre, wie ein Urteil eines komplexen Falles. Es ergibt sich aus der binären Natur des Modells, dass es eher für Ja- oder Nein-Entscheidungen dient, zumal damit keine weitgehende Fallprüfung durchgeführt werden könnte. Allenfalls wäre es möglich das Modell im schweizerischen Recht anzuwenden. So könnte die Methode unter Umständen durch gewisse Abänderungen oder Erweiterungen als Hilfsmittel bei der Verletzung von Grundrechten der Schweizerischen Bundesverfassung dienen. Dazu könnten die vom Schweizerischen Bundesgericht auf dessen Webseite seit 1954 veröffentlichten Fälle verwendet werden.³⁰ Da der von MEDVEDEVA, VOLS und WIELING geschriebene Code offengelegt wurde, könnte dieser übernommen und für die schweizerischen Gerichtsurteile repliziert werden.

²⁴ Vgl. MEDVEDEVA/VOLS/WIELING, S. 254, Fig. 3.

²⁵ Vgl. MIKOLOV et al., S. 3111 ff.

²⁶ Vgl. MEDVEDEVA/VOLS/WIELING, S. 262 f.

²⁷ Vgl. LECUN/BENGIO/HINTON, S. 436 ff.

²⁸ Vgl. JURI Says.

²⁹ Vgl. MEDVEDEVA/VOLS/WIELING, S. 245.

³⁰ Vgl. BGE Rechtsprechung.

5. Fazit

[16] In der vorliegenden Arbeit wurde anhand der Studie von MEDVEDEVA, VOLS und WIELING untersucht, inwiefern Vorhersagen von Machine Learning Methoden sinnvoll erscheinen, um Gerichtsentscheide des EGMR vorherzusagen. Dazu dienten drei verschiedene Experimente, welche mittels *Natural Language Processing* und der *Klassifikationsmethode* basierend auf vergangenen Fällen des EGMR *Verstösse* oder *keine Verstösse* von neun verschiedenen Gesetzesartikel der EMRK vorhersagten. Es konnte gezeigt werden, dass durch die Umwandlung relevanter Textelemente von Gerichtsurteilen in *n-grams* deren Verfahrensausgänge mit einer Genauigkeit von 75 % vorhergesagt werden konnten. Die korrekte Schätzung von künftigen Urteilen aufgrund von vergangenen Gerichtsentscheiden wies jedoch einige Schwierigkeiten auf, während alleine basierend auf den Nachnamen der an den Urteilen beteiligten Richtern je nach Gesetzesartikel präzise Vorhersagen erfolgten. Schliesslich wurden die Vor- und Nachteile des *Natural Language Processing* in Kombination mit *Klassifikation* gegeneinander abgewogen und die Zukunftschancen einer solchen Methode für schweizerische Gerichtsentscheide aufgezeigt.

JANICE DEBORAH KOWALSKI, B.A. ist Studentin der Rechtswissenschaften (MLaw) an der Universität St. Gallen (HSG).