



miracum

Medical Informatics in Research and Care in University Medicine

FROM DATA TO KNOWLEDGE – STRATIFIZIERTE SUBGRUPPEN FÜR DIE ENTWICKLUNG VON PRÄDIKTIONSMODELLEN

Erfahrungen, Erkenntnisse und Berichte 2018 bis 2020

USE CASE

2





Daniela Zöllner



Harald Binder

Fotos: Sommer

Vorwort

Das MIRACUM Konsortium hat für die Aufbau- und Vernetzungsphase der Medizininformatik-Initiative (MI) drei Use Cases definiert und darüber Anwendungsszenarien beschrieben, welche auf den Basiskomponenten der MIRACUM Datenintegrationszentren (DIZ) aufsetzen und so deren Nutzen durch ergänzende innovative IT-Lösungen belegen.

In Use Case 2 „**From Data to Knowledge – Stratifizierte Subgruppen für die Entwicklung von Prädiktionsmodellen**“ werden mit Techniken des maschinellen Lernens (insbesondere mit Deep Learning) valide Vorhersagemodelle auf Grundlage einer großen Fülle an Daten entwickelt. Der schrittweise inhaltliche Ausbau der DIZ an den MIRACUM Standorten bietet dabei eine solide Datenbasis, um Patientenkohorten anhand klinischer Parameter, Biomarker und molekularer/genomischer Untersuchungen in Subgruppen zu stratifizieren. Das Konsortium hat sich auch zur Aufgabe gemacht, entstehende Vorher-

sagemodelle schnellstmöglich mittels Smart-Apps in den Klinikalltag zurückzuspielen, um Ärzt:innen in ihren diagnostischen und therapeutischen Entscheidungen zu unterstützen. Der klinische Fokus liegt hierbei zunächst auf medizinische Fragestellungen aus dem Bereich Asthma/COPD und auf Hirntumoren. In dieser Broschüre stellen wir unsere Ziele, bisher etablierte Konzepte und die derzeitigen Umsetzungsergebnisse vor. Auch bedanken wir uns hiermit bei unserem Use Case 2-Team für die letzten dreieinhalb Jahre sehr engagierter und fruchtvoller Zusammenarbeit.

Daniela Zöllner

Harald Binder

Use Case 2 Koordinatoren

Hans-Ulrich Prokosch

Konsortialleiter MIRACUM

Herausgeber

Steering Board des MIRACUM Konsortiums

Artdirection

W.A.S.

Illustrationen

Nina Eggemann

Druck

Pinguin Druck GmbH

Printed in Germany

Nachdruck, auch auszugsweise, Aufnahme in Onlinedienste und Internet sowie

Vervielfältigungen auf Datenträgern wie CD-ROM, DVD-ROM etc. nur nach vorheriger schriftlicher Zustimmung der Herausgeber. Germany 2021

Verantwortliche

Prof. Dr. Harald Binder

Er hat in Regensburg und an der University of California, Irvine, Psychologie und Mathematical Behavioral Sciences studiert und promovierte 2006 an der LMU München. Nach seiner Postdoc-Zeit in Freiburg und einer Professur an der Universitätsmedizin Mainz trat er 2017 die Professur für Medizinische Biometrie und Statistik am Uniklinikum Freiburg an, verbunden mit der Leitung des gleichnamigen Instituts. Sein Forschungsschwerpunkt ist die Verknüpfung statistischer Techniken mit Ansätzen des Deep Learning für molekulare und klinische Daten.



Prof. Dr. Till Acker

Er hat Medizin in Freiburg, London, San Diego und Kapstadt studiert. Seit 2008 ist er Direktor des Instituts für Neuropathologie an der Justus-Liebig-Universität Gießen, seit 2015 medizinischer Forschungsdekan und seit 2019 Vorsitzender der Deutschen Gesellschaft für Neuropathologie und Neuroanatomie (DGNN). Sein Forschungsschwerpunkt liegt auf der molekularen und funktionellen Kartierung der Rolle des Mikromilieus in der Regulation von „Hallmarks of Cancer“. Klinisch ist seine Gruppe an der Aufdeckung morphomolekularer Biomarker durch angewandte KI für die Differentialdiagnose und Prognose bei Hirntumoren interessiert.



Dr. Daniela Zöller

Sie hat an der Hochschule Koblenz Biomathematik studiert und 2017 an der Universitätsmedizin Mainz im Bereich Biostatistik promoviert. Nach einer 3-jährigen Postdoc-Zeit im Bereich Knowledge Discovery and Synthesis des Instituts Medizinische Biometrie und Statistik des Universitätsklinikums Freiburg, hat sie 2021 die Leitung des Bereichs Medical Data Science übernommen. Ihr Forschungsschwerpunkt liegt im Bereich der Statistischen Methodenentwicklung für Routine- und Registerdaten und für verteilte Analysen unter Datenschutzbeschränkungen.



Prof. Dr. Harald Renz

Seit 1999 ist er Direktor des Instituts für Laboratoriumsmedizin am Universitätsklinikum Gießen und Marburg (UKGM) und Professor der Philipps-Universität Marburg. Seit 2010 ist Prof. Renz stellvertretender Sprecher des UKGM Lung Center. Sein Forschungsinteresse gilt der Pathogenese von Allergien und Asthma. In den Jahren 2012/13 war er Gastprofessor und Fulbright-Stipendiat an der Harvard Medical School in Boston. Seit 2015 ist er CMO des UKGM, Standort Marburg. Von 2010 bis 2016 war er Präsident der Deutschen Gesellschaft für Allergologie und Klinische Immunologie (DGAKI) und seit 2018 ist er Vizepräsident der Deutschen Gesellschaft für Laboratoriumsmedizin. Derzeit hält er Gastprofessuren in Moskau, Russland und Moshi, Tansania.



Team

Christian Bruns

Roland Buhl

Alejandro Cosa Linan

Hildegard Dohmen

Frederike Euchner

Kiana Farhadyar

Patrick Fischer

Denis Gebele

Timm Greulich

Julian Gründner

Dennis Hasenpflug

Christian Haverkamp

Petar Horki

Nattika Jugkaeo

Stefanie Korn

Stefan Lenz

Anke Lux

Michael Maxheim

Attila Nemeth

Michael Neumaier

Alan Race

Stephan Ringshandl

Jannik Schaaf

Jan Scheer

Stefanie Schild

Sebastian Schindler

Bernd Schmeck

Amir Tabatabaei

Dativa Tibyampansha

Dennis Toddenroth

Frederik Trinkmann

Abishaa Vengadeswaran

Johannes Wolf

Jochen Zohner



Vorhersagen entwickeln und klinische Prozesse revolutionieren

Das Gesundheitswesen produziert gigantische Datenmengen mit vielen bislang unbekanntem Informationen. Statistische Werkzeuge, wie z.B. Deep Learning, können aus diesen ungehobenen Schätzen potenziell präzise Vorhersagen über Krankheitsverläufe oder Therapieoptionen in praxisrelevanten Situationen ermitteln. Ziel des Use Case 2 ist es, relevante Datenmuster zu identifizieren und zu validen Vorhersagemodellen zu kombinieren.

Aufgrund individueller Unterschiede, beispielsweise in der Ansprechrate auf Medikamente, ist oftmals die optimale Therapiekombination noch unbekannt und zahlreiche Krankheitsbilder in Deutschland werden noch immer nach Trial-and-Error therapiert. Ein System, das für alle Beteiligten gleichermaßen unbefriedigend ist: Patient:innen haben oftmals mit unerwünschten Nebenwirkungen von Medikamenten zu kämpfen, die letztendlich nur suboptimal wirken; Ärzt:innen wollen den Patient:innen schnelle und wirkungsvolle Therapien anbieten; und auch für die Pharmaindustrie wird es zum Problem, sich den Vorwürfen ausgesetzt zu sehen, unwirksame Mittel zu vertreiben. Dazu kommt, dass dieses Vorgehen tatsächlich auch volkswirtschaftlich eigentlich nicht zu verantworten ist, da die

Ressourcen des Gesundheitswesens nicht optimal genutzt werden. Aus diesem Grund ist es wünschenswert, Patient:innen, anhand von Biomarkern und -signaturen in therapeutisch relevante Subgruppen (Endotypen) einteilen zu können, um sie effizienter und individueller zu behandeln.

Noch zu viel Theorie für die Praxis

Im klinischen Kontext ist eine immer größere Menge an Patientendaten elektronisch verfügbar. Sowohl klinische als auch Labor-daten, Röntgenaufnahmen und individuelle Krankheitsbilder sind dabei oft sogar im Zeitverlauf vorhanden, so dass prinzipiell hoch-informative Profile vorliegen. Um jedoch umsetzbares Wissen zu generieren, müssen aus diesen Daten mit Verfahren der Statistik und der künstlichen Intelligenz Muster iden-

Hirntumore im Visier

Hirntumore werden als besonders belastend empfunden und gehen mit einer hohen Morbidität und Mortalität sowie hohen sozioökonomischen Kosten einher. Bei Kindern und Jugendlichen sind Hirntumore der zweithäufigste Krebstyp. Bei Erwachsenen ist der häufigste Hirntumor das Glioblastom mit einer infausten Prognose und medianen Überlebenszeit von 12-15 Monaten und einer der schlechtesten Fünf-Jahres-Überlebensraten unter allen Krebserkrankungen. Bis heute ist die korrekte Diagnose von Hirntumoren herausfordernd, auch bei erfahrenen Neuropathologen, was die Notwendigkeit herausstreicht, unsere Fähigkeit zur Diagnose und adäquaten Therapie von Hirntumoren zu verbessern.

tifiziert werden, welche für die Behandlung von Patient:innen relevant sind.

Aus solchen Mustern können Diagnose- und Vorhersagemodelle entwickelt werden, die zurück in den Routineeinsatz übertragen werden müssen. Trotz der Fortschritte bei der Entwicklung derartiger Modelle in den letzten Jahren, ist es immer noch eine Herausforderung, auch tatsächlich den Kreis bis in die klinische Routine zu schließen. So markiert in vielen Fällen schon die Bewertung der prinzipiellen Einsetzbarkeit von Vorhersagemodellen bereits den Abschluss solcher Forschungsprojekte. Nur selten wird versucht, die entwickelten Modelle auch tatsächlich in der Krankenversorgung zum Einsatz zu bringen. Um das Potenzial wirklich auszuschöpfen, ist es wichtig, diesen Schritt zu gehen und Algorithmen und Tools für die tägliche Entscheidungsunterstützung einzusetzen und zu verbreiten.

Im Use Case 2 wird MIRACUM deshalb nicht nur Vorhersagemodelle entwickeln, sondern diese in die klinischen Versorgungsprozesse integrieren. Ziel ist es, für mindestens zwei große Krankheitsbilder zu demonstrieren, wie Vorhersagemodelle entwickelt, trainiert und evaluiert werden können, und wie diese in innovative IT-Lösungen überführt werden können, die Ärzt:innen bei konkreten Entscheidungen unterstützen.

MIRACUM hat mit Asthma/COPD (Chronisch obstruktive Lungenerkrankung) und Neuroonkologie den thematischen Schwerpunkt auf zwei medizinische Bereiche gelegt, die gute Beispiele für die Integration hetero-

» Wir möchten im Use Case 2 für mindestens zwei große Krankheitsbilder Vorhersagemodelle mit Deep Learning entwickeln, trainieren und evaluieren. Und wir möchten zeigen, wie unsere Ergebnisse als innovative IT-Lösungen die Mediziner bei konkreten Entscheidungen unterstützen. «

Harald Binder

gener Datenquellen zur Identifizierung prognostisch relevanter Subgruppen darstellen. Obwohl dies zwei unterschiedliche medizinische Spezialgebiete sind, erfordern sie einen gemeinsamen methodischen Kern, welcher in den Datenintegrationszentren (DIZ) der MIRACUM-Standorte etabliert werden wird. Beide medizinischen Fragen profitieren von der Kooperation mehrerer Kliniken. Für beide Indikationen gilt, je größer die Fallzahlen, desto besser lassen sich klinisch relevante Muster identifizieren und bewerten. So ist z.B. die Existenz bestimmter Patientenuntergruppen oft bereits aus Forschungsdatensätzen bekannt, doch noch ist unklar, wie solche Untergruppen basierend auf Routinedaten erkannt werden können und ob Patient:innen zuverlässig als Mitglied bestimmter Untergruppen identifiziert werden können.

Asthma & COPD

Asthma und COPD gehören zu den häufigsten chronisch-entzündlichen Lungenerkrankungen. Jüngste Forschungsergebnisse weisen auf eine große Heterogenität der

zellulären und molekularen Entzündungssignalwege bei Patientenuntergruppen hin, welche als „Endotypen“ bezeichnet werden. Die Endotypisierung von Asthma- und COPD-Patient:innen entwickelt sich so zu einer vorrangigen Aufgabe, da darauf aufbauende gezielte Therapien die Möglichkeit bieten, strategische Mediatoren in diesen Entzündungsnetzwerken selektiv und spezifisch zu beeinflussen. Eine Präzisionsmedizin bei Asthma und COPD basiert auf der Charakterisierung solcher Endotypen unter Verwendung von Biomarkern. Eine besondere Herausforderung besteht darin, dass eine adäquate Charakterisierung der Patient:innen die Integration longitudinaler Daten erfordert, die verschiedenen Krankheitszuständen entsprechen. Leider ist auch die longitudinale sektorübergreifende Datenspeicherung und Datenintegration ein Gebiet, in dem sich Deutschland bislang nicht als Vorreiter auszeichnete.

Eine weitere Herausforderung in diesem Zusammenhang ist die Entwicklung von Biomarker-Panels, die Endotypen präzise

Künstliche Intelligenz für Prädiktion

Die Erstellung von Prognosen oder prädiktiven Vorhersagen spielt in vielen Bereichen des täglichen Lebens eine Rolle. So benutzen E-Commerce-Anbieter Techniken der künstlichen Intelligenz/ des Deep Learnings um das Kaufverhalten von Kunden auf Basis der Historie besuchter Webseiten vorherzusagen. In ähnlicher Weise können biomedizinische Parameter als „Marker“ bzw. „Signaturen“ verwendet werden, um individuell auf bevorstehende Krankheitsveränderungen hinzuweisen bzw. abhängig vom Krankheitsbild bzw. -verlauf individuell passende und maßgeschneiderte therapeutische Maßnahmen zu ergreifen.

widerspiegeln. Eine wichtige Rolle kommt dabei Eosinophilen im Blut zu, wobei gerade Fragen bei nicht-eosinophilem Asthma und COPD aufgeworfen wurden. Neuere klinische und experimentelle Daten legen nahe, dass diese Patientenpopulation nicht nur größer als erwartet ist, sondern auch innerhalb eines Endotyps eine breite Heterogenität besteht. Was erklärt, dass eine erfolgreiche Therapie dieser Patient:innen von Variablen abhängen kann, die wir bisher nicht immer im Griff haben. Für uns ergibt sich daraus auch gezielte Therapien für diese wichtigen Untergruppen von Asthma-/COPD-Patient:innen mit im Auge zu haben.

Um also klinisch relevante Asthma- und COPD-Endotypen besser zu definieren, haben wir den MIRACUM-DIZ-Datensatz um Elemente von Asthma- und COPD-Patient:innen erweitert, wobei insbesondere Parameter der Lungenfunktionsmessung zu nennen sind. Alle Universitätskliniken, die an MIRACUM teilnehmen, bieten bereits bei Erwachsenen- und/ oder Kinderasthma-/ COPD-Patient:innen eine (spezialisierte) Betreuung an. Insgesamt erwarten wir mehr als 3.500 Fälle pro Jahr. Hier zeigt sich der Vorteil einer großen Zahl von Zentren, um eine Unterklassifizierung von Patientengruppen zu ermöglichen. Die Integration heterogener Datenquellen ermöglicht dabei eine tiefe Immunphänotypisierung.

Als Ergebnis werden wir die erste Plattform entwickeln, die die klinische Beratung von Asthma-/ COPD-Patient:innen basierend auf klinischem Clustering ermöglicht. Dies wird

» Wir wissen, dass das One-size-fits-all-Modell in der Therapie oft nicht funktioniert. Die großen Fallzahlen helfen uns, klinisch relevante Muster besser identifizieren und bewerten zu können. «

Harald Renz

nicht nur eine wichtige Basis für die Auswahl und Implementierung von Behandlungsoptionen bieten, sondern langfristig auch zu einer Risikobewertung und Vorhersage der Patient:innen hinsichtlich ihrer klinischen Ursache führen, z.B. eine Entwicklung von Asthma Exazerbation als Haupttreiber der Schwere der Erkrankung.

COPD: Unterschiede zwischen Patient:innen mit und ohne Alpha-1-Antitrypsin-Mangel

Mittlerweile haben sich durch die enge Zusammenarbeit mit Klinikern an den Standorten Freiburg, Marburg und Mannheim mehrere wissenschaftliche Fragestellungen ergeben, die kurz davor sind publiziert zu werden. Ein spezieller Fokus liegt dabei auf COPD Patient:innen mit einer Erbkrankheit, Alpha-1-Antitrypsin-Mangel (AATM), der bewirkt, dass die Leberzellen das Enzym Alpha-1-Antitrypsin fehlerhaft oder in zu geringer Menge bilden oder freisetzen. Dieser Enzymmangel kann dann bereits in jungen Jahren eine COPD auslösen, wodurch sich solche Patient:innen grundlegende von anderen

COPD Patient:innen unterscheiden, welche zumeist älter sind und eine starke Raucherhistorie haben. In der ersten Fragestellung vergleichen wir die Ergebnisse innerhalb des MIRACUM Konsortiums mit denen des COSYCONET Registers. Die Daten des COSYCONET Registers sind im Gegensatz zu den MIRACUM Daten speziell zu Forschungszwecken erhoben worden und bilden dadurch einen Goldstandard. Wir konnten die wichtigsten Ergebnisse reproduzieren: Beispielsweise ist bekannt, dass Patient:innen mit AATM ein erhöhtes Risiko für Lungenemphysemen und Bronchiektasien im Vergleich mit Patient:innen ohne AATM aufweisen, was sich auch in den MIRACUM Analysen gezeigt hat. Umgekehrt konnte wie erwartet kein Unterschied beim Auftreten von Lipoprotein oder Diabetes beobachtet werden. Dies zeigt, dass die Routinedaten grundsätzlich plausibel sind und somit, bei vorsichtiger Interpretation, in der Forschung verwendet werden können. Des Weiteren haben wir erste Hinweise auf das große Potential der Routinedaten zur Unterstützung der Gewinnung neuer medizinischer Erkenntnisse gefunden: Der Anteil

an Pneumokokken- und Pseudomonasinfektionen ist bei Patient:innen mit AATM im Verhältnis höher als bei den Patient:innen ohne AATM. Auch wenn diese Ergebnisse aktuell noch als Work-In-Progress zu sehen sind, könnte diese Beobachtung einen Hinweis auf unterschiedliche Infektionsprofile und damit unterschiedliche klinische Behandlungsfokussierungen liefern.

Asthma & COPD: Der Einfluss von chronischen Lungenerkrankungen auf den Verlauf von hospitalisierten Patient:innen mit COVID-19

Im Zuge der COVID-19 Pandemie haben sich neue Fragestellungen entwickelt, für die Use Case 2 bereits wichtige Vorarbeiten geleistet hatte: Wie wird der Verlauf der COVID-19 Erkrankung im Krankenhaus durch chronische Lungenerkrankungen wie Asthma und COPD beeinflusst? Um dies zu adressieren, wurde der Use Case 2 Datensatz für Asthma und COPD um COVID-19 relevante Parameter wie Beatmungsdauer und Krankenhausmortalität erweitert. Erste vorläufige Ergebnisse zeigen, dass der Anteil an Asthma Patient:innen innerhalb der COVID-19 Kohorte unter der deutschlandweiten Prävalenz liegt, während dies bei COPD Patient:innen nicht der Fall ist.

Hirntumore im Visier

Das zweite medizinische Anwendungsfeld liegt im Bereich der Hirntumore. Sie haben einen direkten und zutiefst beeinträchtigenden Einfluss auf die kognitiven Funktionen,

die geistigen Fähigkeiten und die Persönlichkeit der Patient:innen. Neben diesen hohen Morbiditätsraten sind sie mit einer hohen Mortalität und hohen sozioökonomischen Kosten verbunden. Eine präzise Tumordiagnostik ist daher von zentraler Bedeutung für die richtige Behandlung der Patient:innen.

Bis heute ist die korrekte Diagnose von Hirntumoren herausfordernd, was die Notwendigkeit hervorhebt, unsere Fähigkeit zur Diagnose und adäquaten Therapie von Hirntumoren zu verbessern. Moderne Hochdurchsatzanalysen von großen Hirntumor-Kohorten haben gezeigt, dass DNA-Methylierung als robuste Methode zur Klassifizierung verschiedener Tumoreinheiten und zur Entdeckung neuartiger, molekular unterschiedlicher Tumorsubtypen verwendet werden kann. Hierbei werden pro Hirntumor parallel bis zu 850.000 DNA-Methylierungsstellen untersucht und analysiert. In Zusammenarbeit mit dem Deutschen Krebsforschungszentrum (DKFZ) in Heidelberg kombinieren wir darauf aufbauend DNA-methylierungsbasierte integrative Hirntumordiagnosen mit klinischen und longitudinalen Schlüsselparametern, um den Einfluss der molekularen Hirntumor-Klassifikation auf die Diagnostik und Behandlung von neuroonkologischen Patient:innen zu untersuchen.

Mit 7.000 Neuerkrankungen pro Jahr in Deutschland gehören Hirntumore zu den selteneren Tumorerkrankungen, was die Identifizierung von Subtypen weiter erschwert. Die Größe des MIRACUM Konsortiums ist hierbei für unsere Forschung ein

» Prädiktionsmodelle benötigen strukturierte Daten, die wir innerhalb von MIRACUM liefern können. Dies ist eine essentielle Voraussetzung, um diagnostische und therapeutische Ansätze in der Präzisionsmedizin mit Hilfe von KI optimieren zu können. «

Till Acker

unschätzbaren Vorteil. Zu den wichtigsten Fragen gehört die Untersuchung neuartiger Tumor-Untergruppen und ihrer klinischen Verläufe. Für die klinische Anwendung wird schließlich entscheidend sein, ob methylierungsbasierte Hirntumorklassen Ärzt:innen dabei unterstützen können, Patient:innen, die wahrscheinlich innerhalb der nächsten 12 Monate versterben werden, zuverlässiger zu identifizieren, um sie besser palliativ versorgen zu können.

Durch Analyse dieser hochkomplexen Daten mit Techniken des Deep Learning werden Schlüsselfragen zu den Auswirkungen der verfeinerten, integrativen histomolekularen Diagnostik auf die Therapie und den Outcome von Hirntumorpatient:innen beantwortet. Da auch in anderen medizinischen Bereichen verstärkt molekulare Messungen aus Hochdurchsatzverfahren (omics), wie z.B. zur DNA-Methylierung oder NGS (next generation sequencing) Verfahren, verwendet werden, betrachten wir die beschriebenen Analysen dementsprechend als beispielhaft und zentral für weitere zukünftige Anwendungen.



Internationale Kooperationen auf dem Weg zum verteilten maschinellen Lernen

Valide Prädiktionsmodelle, wie sie in Use Case 2 entwickelt werden sollen, benötigen große Datenmengen, welche oftmals nur durch die Kombination diverser Datenquellen erstellt werden können. Die Herausforderung besteht im Extrahieren relevanter Subgruppen aus heterogenen Datenquellen – dafür kooperiert MIRACUM jetzt auch mit Pionieren aus Newcastle.

Neue Techniken des maschinellen Lernens sollen komplexere Mechanismen und Muster identifizieren, um Empfehlungen für individuellere Behandlungen von Patient:innen geben zu können. Um neue Erkenntnisse zu generieren, sind diese Verfahren aber auf große, freiwillig gespendete Datenmengen (Stichwort „Big Data“) angewiesen. Durch die Erarbeitung eines Datenschutzkonzeptes an den einzelnen Standorten des MIRACUM Konsortiums ist es möglich, die sensiblen Patientendaten unterschiedlicher Standorte mithilfe verteilten Rechnens mit DataSHIELD zu analysieren und Prädiktionsmodelle zu entwickeln.

Die Analyse soll zu den Daten

Der große Datenschatz, der derzeit verteilt auf die zehn Universitätskliniken im MIRACUM Konsortium liegt, soll nicht in einen zentralen Speicher zusammengeführt, sondern verteilt über alle Standorte für Auswertungen genutzt werden. So kann die notwendige kritische Masse an Daten für die Beantwortung bestimmter Fragestellungen zustande kommen. Das Motto und datenschutzrechtliche Grundkonzept des MIRACUM Konsortiums lautet daher „Wir wollen nicht die Daten zur Analyse bringen, sondern die Analysen zu den Daten.“ D.h. Analysen werden dezentral durchgeführt, ohne dass Kliniken die Patientendaten selbst herausgeben müssen. Viele der Algorithmen zur Datenauswertung lassen sich so umformulieren, dass sie auch auf verteilte Daten angewendet werden können.

Zwischen- und Endergebnisse dieser Analysen lassen keine Rückschlüsse auf einzelne

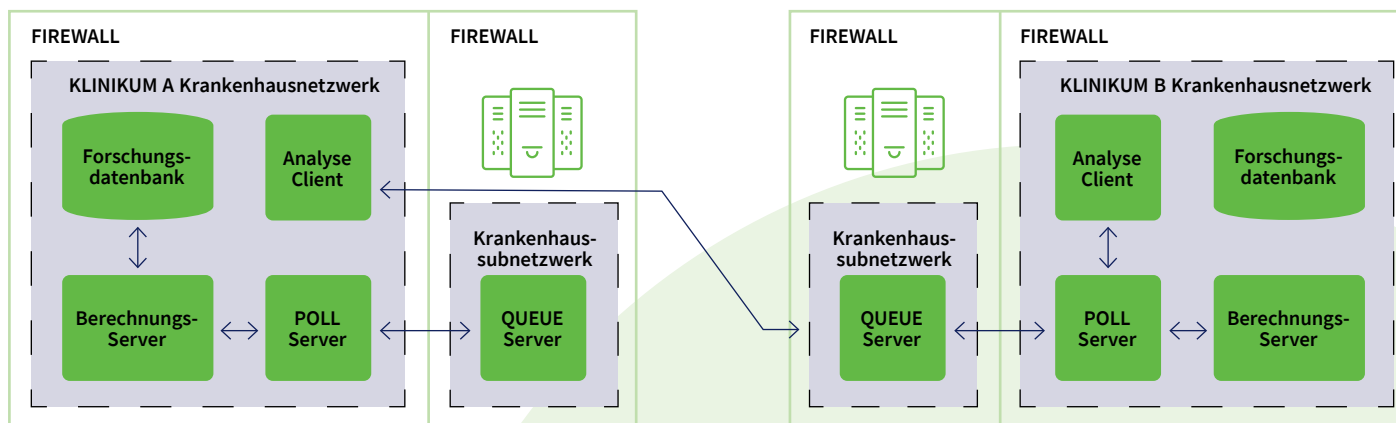
» Eine zentrale Datenbank wäre ein sehr prominenter Angriffspunkt und sollte womöglich schon aus Sicherheitsgründen vermieden werden. «

Patient:innen zu. Dabei werden die teilnehmenden Institutionen (Krankenhäuser) zu einem Analysenetzwerk verbunden, so dass Wissenschaftler:innen über die Plattform Analysen in den teilnehmenden Krankenhäusern durchführen können, ohne Individualdaten der Patient:innen zu sehen.

Das Projekt DataSHIELD

Um solche Analysen durchzuführen, muss eine entsprechend sichere Infrastruktur zwischen den Kliniken aufgebaut werden. Eine Forschergruppe aus England hat auf diesem Gebiet bereits Pionierarbeit geleistet und die Open-Source-Software „DataSHIELD“ entwickelt (Gaye et al. 2014, Budin-Ljøsne et al. 2015). Diese Software bietet bereits verschiedene Verfahren an, die zum statistischen Handwerkszeug gehören. Angefangen bei der Berechnung einfacher Kennzahlen, wie Durchschnittswerten oder Häufigkeiten, bis

VERTRAUENSWÜRDIGE NETZE



Dieser Queue-Poll-Mechanismus erlaubt es Krankenhäusern einfach einem Analysenetzwerk beizutreten, ohne die bestehenden Sicherheitsstandards der Institutionen absenken zu müssen. Auf diese Weise werden im MIRACUM Konsortium mehrere Standorte sicher zu einem vertrauenswürdigen Netzwerk verbunden.

hin zu komplexeren Regressionsmodellen. Zusätzlich zu diesen fertigen Analyseverfahren bietet DataSHIELD aber auch eine flexible Infrastruktur, um neue Arten von Analysen zu entwickeln, die dann auf vernetzte, aber über mehrere Standorte verteilte Datenbestände angewendet werden können. Hier setzen die Use Case 2 Kernentwickler:innen aus Freiburg (Biometrie: Maschinelles Lernen) und Erlangen (Medizininformatik: sichere IT-Infrastruktur) an und bringen ihre Expertise ein.

Erweiterung des DataSHIELD-Werkzeugkastens

Um die Analysen zu den Datenbeständen in einem Netzwerk an die einzelnen Standorte zu bringen, nutzt das Original DataSHIELD-Konzept einen Zugriff auf das jeweilige Klinikumsnetzwerk von außen, um die Analysen

dann im lokalen Netzwerk bereit zu stellen. Dies würde aus Sicherheitsgründen in einem deutschen Universitätsklinikum durch die vor dem Klinikumsnetzwerk eingerichtete Firewall abgewiesen. Dieses Problem wurde vom MIRACUM Use Case 2 Team gelöst, indem eine DataSHIELD Erweiterung konzipiert und implementiert wurde, welche diesen Analysen-Verteilungs-/Zugriffsweg umdreht.

Die mit einem zentralen Client erstellten Analysen werden zunächst außerhalb in einer Warteschlange (Queue) abgelegt. Innerhalb des Netzwerks fragt ein Prozess diese Warteschleife ständig ab, um zu prüfen, ob dort eine neue auszuführende Analyse bereitgestellt wurde. Eine dort abgelegte Analyse wird dann von diesem Prozess mittels eines, über die Firewall erlaubten Prozesses, von außen hinein geholt (Poll).

Vertrauenswürdige Netzwerke generieren

Dieser Queue-Poll-Mechanismus erlaubt es Krankenhäusern, einfach einem Analysenetzwerk beizutreten, ohne die bestehenden Sicherheitsstandards der Institutionen absenken zu müssen. Weiterhin erfasst der Mechanismus jede Anfrage und bietet die Möglichkeit, nur bestimmten Analyse Clients (aus anderen Institutionen) eine Anfrage zu erlauben. Auf diese Weise werden im MIRACUM Konsortium mehrere Standorte sicher zu einem vertrauenswürdigen Netzwerk verbunden.

DataSHIELD-Community

Um diese Erweiterung in die DataSHIELD-Community zu bringen, wurde bereits im Mai 2018 ein erstes Abstimmungstreffen mit dem

Kernentwicklerteam (Prof. Paul Burton, Dr. Olly Butters und Dr. Rebecca Wilson) initiiert. Dieses nahm den Konzeptvorschlag sofort positiv auf, und reagierte mit einer freundlichen Aufnahme in die DataSHIELD-Entwickler:innen-Community. Nach dieser ersten Abstimmung konnten die MIRACUM-Entwicklungen schon im November 2018 beim internationalen DataSHIELD-Workshop in Newcastle (UK) präsentiert werden. Es ergab sich auch ein Einblick in weitere EU-Projekte, die ihrerseits Gesundheitsdaten dezentral analysieren wollen. Europaweit gibt es mehrere Studien, die DataSHIELD einsetzen und mit ihrem Feedback helfen, die Software zu verbessern und auch weitere Gruppen, die neue Funktionalitäten beisteuern. Das von der EU und vom kanadischen Gesundheitsministerium geförderte und 2019 angelaufene



Julian Gründner (M. Sc., Erlangen) und Stefan Lenz (M. Sc., Freiburg) präsentieren ihre DataSHIELD Erweiterungen auf dem DataSHIELD Workshop im November 2018 in Newcastle.

nen EUCAN-Connect-Projekt arbeitet daran, mittels DataSHIELD eine Datenplattform nach den Richtlinien der FAIR-Data-Initiative aufzubauen.

TMF Task Force Verteilte Analysen

Das Thema „Verteilte Analysen“ spielt für die gesamte MII eine große Rolle. 2020 wurden zwei eintägige Workshops zu diesem Thema veranstaltet, bei der die verschiedenen Konsortien der MII bisherige Lösungsansätze präsentiert und Vor- und Nachteile vergleichend diskutiert wurden. Neben dem Ansatz DataSHIELD von MIRACUM sind insbesondere der Personal Health Train und Secure Multi Party Computation zu nennen. Im Nachgang des Workshops wurde von der TMF die Task Force Verteilte Analysen gegründet, in der Use Case 2 durch Harald Binder und Daniela Zöllner aus Freiburg vertreten wird. Hier wurden die Ergebnisse zusammengefasst und eine Stellungnahme zum ak-

tuellen Stand verfasst. DataSHIELD wurde dabei als diejenige Lösung identifiziert, die aktuell bereits für die meisten Einsatzzwecke verwendet werden kann. Aus diesem Grund wird innerhalb der MII nun allgemein empfohlen, DataSHIELD an allen Standorten zur Verfügung zu stellen. Seitdem haben sich bereits einzelne Standorte anderer Konsortien an unser Use Case 2 Team gewandt, um die von MIRACUM entwickelte Queue-Poll-Lösung zu nutzen. Allerdings ist ein weiteres Ergebnis der Task Force, dass DataSHIELD in einigen Bereichen aktuell nicht einsatzbar ist, wie z.B. beim Record-Linkage oder bei der Zusammenführung sehr kleiner Fallzahlen. MIRACUM hat daraufhin Kontakt mit dem Kernentwicklerteam in England aufgenommen, das großes Interesse an einer Kooperation mit den anderen Ansätzen bekundet hat. Es gibt also Bewegung rund um DataSHIELD, zu der auch die deutsche Entwicklungsinitiative durch MIRACUM-Elan beigetragen hat.

Fokus: chronische Lungenerkrankungen und Hirntumore

Gemäß dem Motto „Learning by doing“ wird der Fokus zuerst auf die zwei bereits vorgestellten Erkrankungsgruppen Asthma/COPD und Hirntumore gelegt, um die neuen Techniken zu entwickeln und zu erproben. Diese beiden Projekte stellen die Beteiligten vor verschiedene Herausforderungen.

Im Bezug auf Daten von Asthma/COPD-Patient:innen hat man es vor allem mit Daten aus dem klinischen Routinebetrieb zu tun. Hier ist eine enge Zusammenarbeit zwischen Klinikern, Biostatistikern und Medizininformatikern notwendig, um relevante Fragestellungen, dafür geeignete Auswertungsverfahren und Daten zu identifizieren und aus den Krankenhausinformationssystemen so zu extrahieren. Mittlerweile haben sich durch die enge Zusammenarbeit mit Klinikern an den Standorten Freiburg, Marburg und Mannheim mehrere wissenschaftliche Fragestellungen ergeben, die kurz davor sind publiziert zu werden.

Freiburger Biometrie entwickelt neues Werkzeug

Für die Patient:innen mit Hirntumoren liegen ebenfalls umfangreiche genetische Daten vor. Dabei gilt es, aus dieser unübersichtlichen Menge genetischer Merkmale Hinweise auf diejenigen Eigenschaften herauszufinden, die einen Einfluss auf die Schwere der Erkrankung bzw. die Heilungschancen haben könnten. Diese Hinweise können dann in Laborexperimenten weiter untersucht werden, um die

genauen Wirkmechanismen zu erforschen. Dafür entwickelte die Freiburger Biometrie ein neues Werkzeug zum Finden von Einflussfaktoren für DataSHIELD, das seit 2019 an Hirntumordaten erprobt werden kann.

Diese selbst gestellten Herausforderungen dienen auch der Motivation, um eine nachhaltige und flexible Dateninfrastruktur aufzubauen, damit Forschung und Patient:innen von dem verstreuten Datenschatz profitieren. So wird darauf hin gearbeitet, Werkzeuge bereit zu stellen, damit Ärzt:innen und Biolog:innen Krankheiten genauer verstehen können und individuellere Behandlungen ermöglicht werden, ohne dass dabei der Schutz der Gesundheitsdaten vernachlässigt wird.

Foto: iStock/demiro, Wilson

DataShield Literatur

Budin-Ljøsne I, Burton P, Isaeva J, Gaye A, Turner A, Murtagh MJ, et al.: DataSHIELD: an ethically robust solution to multiple-site individual-level data analysis. *Public Health Genomics*. 2015;18(2):87-96.
Gaye A, Marcon Y, Isaeva J, LaFlamme P, Turner A, Jones EM, Minion J et al.: DataSHIELD: taking the analysis to the data, not the data to the analysis. *Int J Epidemiol*. 2014 Dec;43(6):1929-44.

MIRACUM ist eine riesige Chance für Diagnostik und Therapie

Nach drei Jahren MII, in der zumeist infrastrukturelle Vorarbeiten geleistet wurden, kommt der praktische Nutzen dieses Großprojektes für die Mediziner langsam in der Klinik an. PD Dr. Sebastian Fährndrich, Leitender Oberarzt in der Pneumologie des Universitätsklinikums Freiburg, berichtet aus ärztlicher Sicht.

INTERVIEW MIT PD Dr. Sebastian Fährndrich, Universitätsklinikum Freiburg

Welche Daten interessieren Sie besonders, die durch die MII in greifbare Nähe gekommen sind?

Oh, da muss ich nicht lange überlegen. Lungenfunktionsdaten, Diagnosen, Nebendiagnosen, Keime und in naher Zukunft das Auslesen von Arztbriefen. Natürlich müssen der Datenschutz mit an Bord und die Ethikkommission sowie das Data Use & Access Committee glücklich sein. Dann können wir daran gehen, die Daten wissenschaftlich zu nutzen, während die Patient:innen gleichzeitig ein gutes Gefühl haben.

Und diese Möglichkeiten eröffnen die Nutzung der Datenintegrationszentren in MIRACUM?

MIRACUM ist eine wirklich tolle Sache für Diagnostik und Therapie. In der Medizin haben wir oft das Problem, dass es für Vieles

zu wenige Daten gibt. Die MIRACUM Arbeiten liefern eine tolle Datenbank, die uns die Möglichkeit eröffnet, retrospektiv enorme Fallzahlen auszuwerten, um mit diesen gewonnenen Erkenntnissen Patient:innen zukünftig besser behandeln zu können.

Sehen Sie, Patient:innen kommen mit einer schweren Vorerkrankung zu uns. COPD, die sogenannte Raucherlunge, ist eine Volkskrankung und betrifft sehr viele Menschen. Seit wenigen Jahren wissen wir, dass die chronische Raucherlunge eine System- und auch Umwelterkrankung ist. Das bedeutet, sie betrifft eben nicht nur die Lunge. Partikel wandern aus der Lunge durch den Körper und lagern sich in Gefäßen und Organen ein, wo sie überall im Körper kleine Entzündungen auslösen. Mit MIRACUM bekommen wir Zugriff auf so viele Patient:innen und Daten rings um die Patient:innen, so dass wir neue



PD Dr. Sebastian Fährndrich:
„Bessere Datenlage verspricht verbesserte Therapien.“

Zusammenhänge erfahren, die für das Risikomanagement, Diagnostik und die Therapie von enormer Bedeutung sind.

Ich stelle es mir schwierig genug vor, sich mit Ärzt:innen der eigenen Station auf die beste Behandlung zu einigen. Wie sieht eine Abstimmung über verschiedene Standorte hinweg aus?

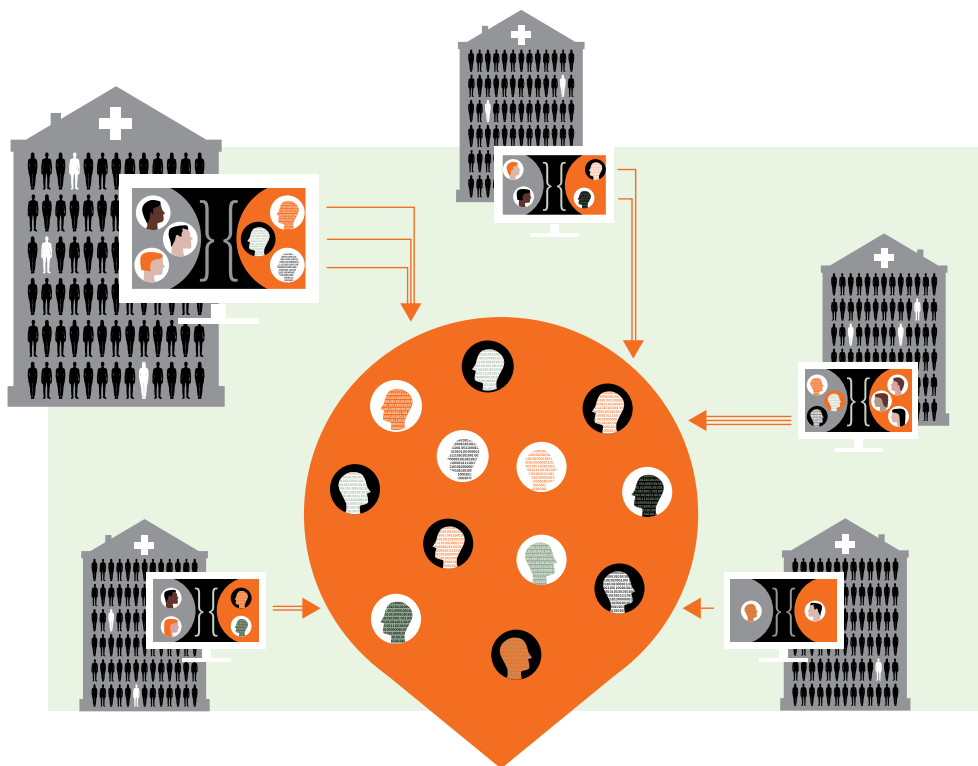
Das ist tatsächlich kein Problem. Wir halten uns an die Leitlinien der Lungenheilkunde, auch speziell für die COPD – zu Schweregrad, Behandlung etc. Behandlungs-Leitlinien sollen mit dem Wissen, welches wir durch die Auswertung der MIRACUM-Daten erwerben, hinterfragt und aktualisiert werden.

Das ermöglicht, die gewonnenen Erkenntnisse auch in die Diagnostik und Behandlung der Patient:innen einfließen zu lassen. Wichtig und neu ist die Zusammenarbeit mit den

Medizininformatiker:innen. Die müssen mitgenommen werden und verstehen, welche klinischen Parameter wichtig sind, so dass sie in das DIZ einfließen müssen. Das ist das Wunderbare an MIRACUM – die Schnittstellen zwischen Ärzt:innen und Informatiker:innen funktionieren – wir können uns aufeinander verlassen, darauf dass wir über die gleichen Datensätze und -inhalte sprechen.

Profitieren die Patient:innen schon heute davon, dass sich im Hintergrund der Kliniken bereits an so vielen unterschiedlichen Stellen abgestimmt wird?

Das ist heute vielleicht doch noch ein wenig zu früh, um Genaueres zu sagen. Konkret profitieren werden sie ab dem Zeitpunkt, wenn neue Erkenntnisse einfließen und die ersten Studien ausgewertet sind. MIRACUM ist eher eine Plattform, die uns ermöglicht, für Risikofaktoren, Diagnostik und Therapie zu schauen, welche Fragestellungen in Zukunft beantwortet werden müssen, um Patient:innen wirkungsvoll auch präventiv zu schützen.



Die virtuellen Patient:innen

Rigide Datenschutzbestimmungen machen in Deutschland der Forschungslandschaft oftmals schwer zu schaffen. Um dennoch mit eigenen Daten forschen zu können, prüft das Team des Use Case 2, ob synthetische Daten Teil der Lösung werden können.

Der Schutz von Gesundheitsdaten ist wichtig, da sie personenbezogene Daten beinhalten, die sowohl Patient:innen als auch deren Familienangehörige betreffen. Datenschutzbestimmungen bilden an-

dererseits jedoch ein kritisches Hindernis für den Zugang zu Patientendaten, die zur Verbesserung der Gesundheitsversorgung genutzt werden könnten. Dieses Problem adressieren wir auch im Use Case 2 des

MIRACUM Konsortiums, dessen Ziel es ja ist, personalisierte Vorhersagemodelle für die medizinische Behandlung zu erstellen. Dazu sollen die Daten mehrerer MIRACUM Standorte gemeinsam analysiert werden. Aktuell ist es allerdings noch so, dass diese aus Datenschutzgründen nicht zusammengeführt werden dürfen. Im Rahmen der technischen Lösungen, die wir zum Umgang mit dieser Problematik entwickeln, wollen wir zunächst Daten chronischer Lungenerkrankungen (Asthma und COPD) analysieren und damit im MIRACUM Projekt wissenschaftliche Fragestellungen auf Daten von Patient:innen mit Hirntumoren untersuchen.

Synthetische Daten sind eine Lösung, die zur Überwindung von Datenschutzhürden genutzt werden können. Synthetische Daten sind Daten, die keine echten Daten enthalten, sondern lediglich allgemeine Merkmale und statistische Beziehungen echter Daten nachbilden. Für die Datennutzung in der Forschung bedeutet dies, dass pro Standort virtuelle Patientendaten erstellt werden, die nicht an die Daten einzelner Patient:innen gebunden sind. Solche Daten können dann gemeinsam genutzt werden und ermöglichen den Einsatz statistischer Standardanalysen, oder auch Techniken der künstlichen Intelligenz.

Generative Modelle

Für die Erzeugung synthetischer Daten aus realen Daten ist eine Art statistischer oder maschineller Lernansatz erforderlich. Konkret verwenden wir generative Modelle, welche die systematische und zufällige Va-

riabilität der Originaldaten abbilden, indem ein Modell der statistischen Verteilung der Daten erstellt wird. Tiefe Techniken des Deep Learnings können sich sehr flexibel unterschiedlichen Datenverteilungen annähern. In MIRACUM verwenden wir unter anderem einen der populärsten generativen Ansätze, sogenannte Variational Autoencoder (VAE).

Wie funktioniert ein Variational Autoencoder?

Die Funktionsweise des VAE kann gut an einem Beispiel erklärt werden: Wir wollen wissen, „wie krank“ ein:e Patient:in ist. Im Idealfall haben wir für alle Patient:innen ein oder zwei abstrakte Werte, die den Grad der Erkrankung repräsentieren. Damit könnten wir Daten für virtuelle Patient:innen generieren, indem wir plausible Werte für eine oder zwei Dimensionen auswählen.

In der realen Welt werden Krankheiten jedoch durch viele verschiedene beobachtete Symptommuster gekennzeichnet. Deshalb wollen wir von Mustern, die in realen Messungen, Beobachtungen oder Symptomen zu sehen sind, auf eine kleine Anzahl dieser niedrigdimensionalen Werte abbilden. Diese latenten Merkmale werden im Grunde nie direkt beobachtet, also treffen wir einige plausible Annahmen, um ein Muster zu finden. Wir nehmen zum Beispiel an, dass wir, wenn wir den Krankheitswert aller Individuen in einem Häufigkeitsdiagramm aufzeichnen, eine Normalverteilung oder Gaußsche Glockenkurve haben. Gemäß dieser Annahme können wir dann aus den Daten die Trans-

formation lernen, die Messungen, Beobachtungen und Symptome in die latenten Krankheitswerte transformiert.

Für den nächsten Schritt stelle man sich vor, dass ein:e Mediziner:in Messungen, Beobachtungen und Symptome von Patient:innen mit bestimmten Erkrankungen mit anderen Wissenschaftler:innen diskutieren möchte, natürlich ohne den Datenschutz zu verletzen. Mit den latenten Krankheitswerten für verschiedene Individuen ist es möglich, eine weitere Transformation vom Krankheitswert zurück zu passenden Messungen, Symptomen oder Diagnosen zu lernen. Damit können schließlich plausibel gewählte latente Krankheitswerte virtueller Patient:innen zurücktransformiert werden, so dass sie wie plausible Patientendaten aussehen.

Herausforderungen bei der Verwendung Variational Autoencoder für klinische Daten

Der VAE hat eine eingebaute Normalverteilungsannahme. Wir wollen VAEs jedoch auch für Daten verwenden, bei denen die ursprünglichen Messungen eine ganz andere Form haben könnten. Es gibt zum Beispiel einige Laborwerte, welche ihr Minimum bei Null haben, im Durchschnitt um 50 liegen und sich aber bis zu Werten im Tausenderbereich erstrecken können, sodass sie nicht einer Normalverteilung folgen. Ein anderes Beispiel sind Messungen, die sich zwischen zwei Gruppen von Personen unterscheiden können, wie es z. B. in Daten zur Lungenfunktion der Fall ist. Ein typischer Messwert in diesen

Daten ist die Vitalkapazität, welche bei Frauen und Männern aufgrund der verschiedenen Körpergrößen unterschiedlich ist. Wenn man nun die Verteilung aufzeichnet, ist sie daher nicht mehr glockenförmig, sondern eher eine Kombination aus zwei Glockenkurven, welche durch eine Kerbe getrennt wären. Da wir realistische virtuelle Patient:innen haben wollen, müssen wir in der Lage sein, plausible Werte für derartige Verteilungen zu erzeugen. Wir müssen in den erwähnten Labordaten die kleineren Werte für Frauen im Vergleich zu Männern in der Vitalkapazität in unseren virtuellen Patient:innen wiederfinden, um sie für weitere Forschungen nutzbar zu machen. Leider haben wir festgestellt, dass VAEs dazu neigen, Daten zu erzeugen, die einer Glockenkurve folgen, auch wenn die Eingabe anders aussieht. Daher mussten wir eine neue Lösung entwickeln, welche diese Herausforderung angeht.

Unsere Methode: PTVA

Das Wichtige an den neuen Methoden zur synthetischen Datenerzeugung ist, dass sie nicht auf eine einzige Art von Daten zugeschnitten sein sollten. Es ist zum Beispiel keine gute Lösung, wenn wir nun einige andere Annahmen treffen, welche wiederum die Ergebnisse bei tatsächlich glockenförmigen Daten verschlechtern würden. Bei der Methode, die wir entwickeln, ändern wir also nicht die Annahmen, sondern versuchen, Transformationen zu finden, die schwierige Datenverteilungen bei Bedarf in glockenförmige Verteilungen umwandeln können.

Um die jeweils beste Transformation für die Daten zu finden, benutzen wir Maßzahlen, die anzeigen, ob deren Verteilung glockenförmig ist oder ob es sich um eine andere Verteilung handelt. Nach diesen Maßzahlen kann man die besten Parameter für die Transformation in eine glockenförmige Verteilung finden. Durch die Anwendung dieser Transformationen, kombiniert mit der Verwendung des VAEs mit Normalverteilungsannahme und der Verwendung von Rücktransformationen erhalten wir nun synthetische Daten für anspruchsvollere Arten von Verteilungen. Auf diese Weise können wir realistische Daten virtueller Patient:innen für viele verschiedene Anwendungen erstellen. Diese Methode führen wir unter dem Namen Pre-transformation Variational Autoencoder (PTVA).

Die virtuellen Patient:innen in DataSHIELD

Die Generierung virtueller Patientendaten muss verteilt über verschiedene MIRACUM-Standorte durchgeführt werden. Dafür und allgemein zur Durchführung gemeinsamer, standortübergreifender Analysen benutzen wir die Software DataSHIELD als Infrastruktur. Diese Software ist unter einer Open-Source-Lizenz veröffentlicht und frei nutzbar. Das MIRACUM-Team steht in regelmäßigem Austausch mit dem Kernentwicklerteam an der Universität in Newcastle.

DataSHIELD bietet bereits verschiedene Verfahren an, die zum statistischen Handwerkszeug gehören, angefangen bei der



Kiana Farhadyar (M. Sc.)

Berechnung einfacher Kennzahlen, wie Durchschnittswerten oder Häufigkeiten, bis hin zu komplexeren Regressionsmodellen. Zusätzlich zu diesen fertigen Analyseverfahren bietet DataSHIELD aber auch eine flexible Infrastruktur, um neue Arten von Analysen zu entwickeln und auf damit vernetzte Daten anzuwenden. Diese können wir auch dafür benutzen, um eigene Ansätze umzusetzen. Im letzten Jahr konnten wir die Generierung synthetischer Daten mittels generativer Modelle bereits an einem anderen Verfahren, den so genannten Deep Boltzmann Machines, in DataSHIELD erproben. Nun werden wir den Werkzeugkasten von DataSHIELD noch um die VAEs erweitern, um qualitativ hochwertige synthetische Daten für diverse Verteilungen von Daten erstellen zu können.



Ausblick

In Use Case 2 konnten wir bisher ein breites Spektrum von praktisch nutzbaren Techniken entwickeln, wie z.B. Deep Learning-Ansätze. Ein wesentliches Erfolgsrezept war dabei der enge Austausch mit Nutzer:innen in Klinik und Forschung. Dies ist auch die Basis, um die Entwicklungen im weiteren Projektverlauf noch enger mit dem klinischen Alltag und medizinischen Fragestellungen zu verzahnen und so das Versprechen von künstlicher Intelligenz in der Medizin einzulösen. —

Zum Weiterlesen ...

Publikation, die im Use Case 2 entstanden sind:

1. Farhadyar K, Bonofiglio F, Zoeller D, Binder H. Adapting deep generative approaches for getting synthetic data with realistic marginal distributions. [export.arXiv.org > stat > arXiv:2105.06907](https://arxiv.org/abs/2105.06907); submitted 14.05.2021
2. Lenz S, Hess M, Binder H. Deep generative models in DataSHIELD. *BMC Med Res Methodol* 21, 64 (2021). <https://doi.org/10.1186/s12874-021-01237-6>
3. Gruendner J, Wolf N, Tögel L, Haller F, Prokosch HU, Christoph J. Integrating Genomics and Clinical Data for Statistical Analysis by Using GEnome MINing (GEMINI) and Fast Healthcare Interoperability Resources (FHIR): System Design and Implementation. *JMIR* 2020;22:e19879. Doi: 10.2196/19879
4. Jaravine V, Balmford J, Metzger P, Boerries M, Binder H, Boeker M. Annotation of Human Exome Gene Variants with Consensus Pathogenicity. *Genes* 11 (9), 1076 (2020). <https://doi.org/10.3390/genes11091076>
5. Gruendner J, Schwachhofer T, Sippl P, Wolf N, Erpenbeck M, Gulden C, Kapsner LA, Zierk J, Mate S, Stürzl M, Croner R, Prokosch HU, Toddenroth D. KETOS: Clinical decision support and machine learning as a service – A training and deployment platform based on Docker, OMOP-CDM, and FHIR Web Services. *PLoS ONE* 14(10): e0223010. <https://doi.org/10.1371/journal.pone.0223010>
6. Gruendner J, Prokosch HU, Schindler S, Lenz S, Binder H. A Queue-Poll Extension and DataSHIELD: Standardised, Monitored, Indirect and Secure Access to Sensitive Data. *Stud Health Technol Inform.* 2019;258:115-119

GEFÖRDERT VOM



Bundesministerium
für Bildung
und Forschung



MITGLIED DER

