

Ninja Marnau, Pascal Berrang, Mathias Humbert

Anonymisierungsverfahren für genetische Daten

Privacy-Enhancing Technologies (PETs) wie Differential Privacy und anderen Anonymisierungsverfahren kommen im Hinblick auf hochsensible Gesundheitsdaten besondere Bedeutung zu. Der vorliegende Beitrag zeigt neue Datenschutzrisiken bei epigenetischen Daten auf, entwickelt und analysiert geeignete Gegenmaßnahmen durch PETs und diskutiert die rechtliche Bewertung von deren Einsatz in der medizinischen Forschung.

Der für viele Big Data Anwendungen immanente Konflikt zwischen einer eventuell sogar gesellschaftlich gewünschten Datenerhebung und –nutzung im öffentlichen Interesse und möglichen hohen Risiken für einzelne Betroffene zeigt sich in besonderer Weise im Bereich der medizinischen Forschung und insbesondere in der Genetik.

¹ Diese Arbeit wurde unterstützt durch das Bundesministerium für Bildung und Forschung (BMBF) im Rahmen der Förderung für das Center for IT-Security, Privacy and Accountability (CISPA) (FKZ: 16KIS0656) und die Deutsche Forschungsgemein-



Ninja Marnau, Ass. iur.

Rechtswissenschaftlerin und Senior Researcher am CISPA – Center for IT-Security, Privacy and Accountability in Saarbrücken (zukünftig CISPA – Helmholtz-Zentrum i.G.). Forschungsschwerpunkte: Privacy-by-Design und IT-Sicherheit.

E-Mail: marnau@cispa.saarland



Pascal Berrang, M.Sc.

Wissenschaftlicher Mitarbeiter am CISPA – Center for IT-Security, Privacy and Accountability in Saarbrücken (zukünftig CISPA – Helmholtz-Zentrum i.G.). Forschungsschwerpunkt: Genetic Privacy

E-Mail: pascal.berrang@cispa.saarland



Mathias Humbert, Ph.D.

Senior Data Scientist am Swiss Data Science Center (SDSC), einem Joint Venture der EPFL und der ETH Zürich. Forschungsschwerpunkt: Genetic Privacy

E-Mail: mathias.humbert@epfl.ch

Seit der ersten Sequenzierung des menschlichen Genoms im Jahr 2001 wurden bereits tausende Genome und über eine Million Genotypen sequenziert. Die medizinische Forschung erlaubt es uns, das Risiko für genetisch (mit)bedingte Krankheiten wie Krebs, kardiovaskuläre und neurodegenerative Erkrankungen immer genauer vorherzusagen und abzuschätzen. Auch die Entwicklung von personalisierter Behandlung und Vorsorge gewinnt mit dem Forschungsfeld der Pharmakogenomik, welche den Einfluss des Genoms auf die Wirkung von Medikamenten betrachtet, immer mehr an Bedeutung.

Die „genetische Revolution“ geht jedoch einher mit hohen Datenschutz-Risiken für die Betroffenen. Die Daten enthalten nicht nur Hinweise auf die Veranlagung zu bestimmten Erkrankungen, sondern stehen auch in direkter Relation zu Phänotyp, Verwandtschaft und ethnischer Abstammung.² All diese Informationen können zu verschiedensten Arten des Datenmissbrauchs und der Diskriminierung von Betroffenen führen, individuell oder kollektiv. Verstärkt wird dieses Risiko noch dadurch, dass genetische Daten sich über die Lebenszeit kaum verändern und gleichzeitig auch Familienmitglieder betreffen.³ Vor dem Hintergrund dieser Bedenken ist es nicht überraschend, dass technischer Datenschutz für genetische Daten bereits Diskussionsgegenstand von zahlreichen Veröffentlichungen ist.⁴

1 Von der Genetik zur Epigenetik

Unser Genom ist nicht das Einzige, was unsere Gesundheit beeinflusst. Umweltfaktoren, wie u. a. Ernährung, Lebensführung und Umweltverschmutzung, spielen bei vielen Erkrankungen eine wichtige Rolle. Die Gebiete der Epigenetik, Transkriptomik und Proteomik versuchen daher diese Lücke zwischen unserer gene-

schaft (DFG) im Rahmen des Sonderforschungsbereichs „Methoden und Instrumente zum Verständnis und zur Kontrolle von Datenschutz“ (SFB 1223), Projekt A5.

² Ayday, De Cristofaro, Hubaux, Tsudik: Whole genome sequencing: Revolutionary medicine or privacy nightmare? *Computer*, 2015, 58–66; Lin, Owen, Altman: Genomic research and human subject privacy. *SCIENCE* 2004, 183.

³ Humbert, Ayday, Hubaux, Telenti: Addressing the concerns of the Lacks family: quantification of kin genomic privacy. In Proceedings of the 2013 ACM SIGSAC CCS 2013, 1141–1152.

⁴ Statt vieler Mittos, Malin, De Cristofaro: Systematizing Genomic Privacy Research – A Critical Analysis, ePrint arXiv:1712.02193 [cs.CR], 2017.

tischen Veranlagung und unserem tatsächlichen Gesundheitszustand zu schließen. Während die DNA-Sequenzierung Auskunft darüber gibt, was eine Zelle potenziell tun kann, erlauben diese Gebiete eine bessere Einsicht darin, was eine Zelle zu einem gewissen Zeitpunkt tatsächlich tut. Eine passende Computer-Analogie für die Epigenetik wäre: Wenn das Genom die Hardware des Körpers ist, dann ist das Epigenom die Software.⁵

Epigenetik bezeichnet ein Gebiet, welches sich hauptsächlich mit zellulären und phänotypischen Veränderungen beschäftigt, die nicht durch Änderungen im Genotyp entstehen. Beispiele für solche externen Faktoren sind die in-utero-Entwicklung, die Entwicklung im Kindesalter, aber auch Umweltchemikalien, Alterung oder Ernährung. Epigenetik kann aber auch die Veränderungen selbst bezeichnen, wie z. B. DNA-Methylierung. MicroRNAs (miRNAs) wiederum sind epigenetisch regulierte Mechanismen, welche in den frühen 1990er Jahren entdeckt wurden. MiRNAs sind kleine, nicht-kodierende RNA-Moleküle, die die Genexpression regulieren. Es wurde gezeigt, dass 60% der menschlichen Gene, die Proteine kodieren, durch miRNAs reguliert werden.⁶ MiRNA-Expressionen sind reelle Zahlen, die in einer zweiseitigen Polymerase-Kettenreaktion (*polymerase chain reaction*, PCR) gemessen werden. Sie messen, wie viele miRNA-Moleküle einer bestimmten miRNA in einer Menge von Zellen oder im Gewebe aktiv sind. Unterschiedliche miRNAs sind in verschiedenen Zellen oder Gewebearten aktiv. Speziell miRNA-Expression – gemessen im Blut von Patienten – wurde bereits erfolgreich mit schwerwiegenden Erkrankungen in Verbindung gebracht. Daher ist es nicht erstaunlich, dass miRNA-Expressionsprofiling eine vielversprechende Technik der Biomedizin ist, um frühere, genauere und minimalinvasive Diagnosen für schwerwiegende Erkrankungen zu erstellen.

2 Neu identifizierte Risiken

Während die Epigenetik im biomedizinischen Forschungsumfeld aktuell immer mehr an Bedeutung gewinnt, wurden Risiken im Hinblick auf die Privatsphäre von Probanden dort bisher als weniger gravierend als in der klassischen Genetik erachtet.

Im Gegensatz zur DNA, die meist über die Zeit unverändert bleibt, glaubten viele Biomediziner, dass sich miRNA-Expressionsprofile gerade durch den Einfluss von Umweltfaktoren über die Zeit genügend verändern, um es unmöglich zu machen, zwei zu unterschiedlichen Zeiten erstellte Expressionsprofile demselben Betroffenen zuzuordnen.

Diese Annahme galt, obwohl epigenetische Daten nicht weniger sensible Informationen als das Genom selbst enthalten. So wurde bereits eine Vielzahl an schwerwiegenden Erkrankungen (z. B. Krebs, Diabetes und Alzheimer⁷) identifiziert, welche direkt mit epigenetischen Veränderungen zusammenhängen. Eine Studie fand sogar heraus, dass epigenetische Veränderungen potenziell einen Zusammenhang mit der sexuellen Orientierung ha-

ben könnten.⁸ Und während genetische Daten hauptsächlich das Risiko einer Erkrankung abschätzbar machen, geben die epigenetischen Daten einen direkten Einblick, ob der Patient eine Erkrankung aktuell in sich trägt.

Dies nahm unser Forschungsinstitut zum Anlass, die Verknüpfbarkeit von miRNA-Expressionsprofilen und somit die Identifizierbarkeit von Betroffenen über längere Zeit zu prüfen.⁹ Unsere Forschungsfrage war, ob es einem Angreifer möglich ist, zwei zu unterschiedlichen Zeitpunkten erstellte miRNA-Expressionsprofile miteinander zu verknüpfen. Auf diese Weise könnte ein Angreifer zum Beispiel eine ihm bekannte Person anhand ihrer miRNA-Daten in den Daten einer Alzheimer-Studie aussondern.

Dieses Risiko ist kein bloßes akademisches Gedankenspiel. Solche miRNA-Expressionsprofile können unter anderem aus öffentlich verfügbaren Forschungsdatenbanken stammen¹⁰ oder aus einem *Data Breach* beim Verantwortlichen. Einer Studie des Ponemon Institute zufolge haben zwischen 2014 und 2016 fast 90% der befragten Gesundheitsdienstleister eine Verletzung ihrer Datensicherheit festgestellt und 56% glauben nicht, adäquate Ressourcen zu haben, um diese Vorfälle zu bekämpfen.¹¹ Angriffe gegen Gesundheitsdienstleister können große Mengen an Patienten betreffen oder zielgerichtet auf bekannte Opfer sein.¹² Kriminell erlangte Gesundheitsdaten sind immer häufiger auch auf dem Schwarzmarkt erhältlich.

Mit der Unterstützung unserer Kollegen aus der Medizin und Bioinformatik konnten wir mehrere hundert große, diverse und über längere Zeiträume erhobene miRNA-Expressionsprofile analysieren. Die Profile enthielten Daten von gesunden Probanden sowie Teilnehmern mit 19 verschiedenen genetisch bedingten Erkrankungen. Unsere Analyse umfasste zwei Arten von Verknüpfungsangriffen: Identifikation und Matching. Während die Identifikation darauf abzielt, ein einzelnes miRNA-Expressionsprofil in einer Datenbank von Profilen, z. B. eines späteren Zeitpunktes, wiederzufinden, gleicht der Matching-Angriff zwei Datenbanken miteinander ab, um Übereinstimmungen und verkettbare Profile zu entdecken.

Es gelang uns erfolgreich zu beweisen, dass miRNA-Expressionsprofile entgegen der häufigen Annahme durchaus auch über längere Zeiträume verkettbar sind.¹³ Da noch nicht viele Langzeitdaten zu miRNA existieren, erstreckten sich unsere Forschungsdaten über ein Maximum von 18 Monaten. Potentiell könnte eine

8 Ngun et al.: A novel predictive model of sexual orientation using epigenetic markers. American Society of Human Genetics 2015 Annual Meeting, 2015.

9 Die im Folgenden dargestellten Forschungsergebnisse wurden in Zusammenarbeit mit Anne Hecksteden (Institut für Sport- und Präventivmedizin, Universität des Saarlandes), Andreas Keller (Klinische Bioinformatik, Universität des Saarlandes) und Tim Meyer (Institut für Sport- und Präventivmedizin, Universität des Saarlandes) sowie Michael Backes (CISPA – Helmholtz-Zentrum i. G.) erzielt. Dieser Artikel greift die Ergebnisse des gemeinsamen Konferenzbeitrags auf der USENIX Security 2016 „Privacy in epigenetics: Temporal linkability of microRNA expression profiles“ auf.

10 Prominente Beispiele sind die Gene Expression Omnibus (GEO) <https://www.ncbi.nlm.nih.gov/geo/> und ArrayExpress (<https://www.ebi.ac.uk/arrayexpress/>) Datenbanken.

11 Ponemon Institute: Sixth Annual Benchmark Study on Privacy & Security of Healthcare Data, Mai 2016.

12 Süddeutsche.de, Krankenakte von Michael Schumacher – Gestohlene Details, 24. Juni 2014, <http://www.sueddeutsche.de/panorama/krankenakte-von-michael-schumacher-gestohlene-details-1.2013913>.

13 Eine Verkettbarkeit oder Identifizierung gelang je nach Art des Angriffs für bis zu 80% der Profile (bei blutbasierten Profilen) oder bis zu 40% der Profile (bei plasmabasierten Profilen).

5 Cloud: Why your DNA isn't your destiny. *Time*, 2010.

6 Friedman, Farh, Burge, Bartel: Most mammalian mRNAs are conserved targets of microRNAs. *Genome research*, 2009, 19(1):92–105.

7 Feinberg, Fallin: Epigenetics at the crossroads of genes and the environment. *JAMA*, 2015, 314:1129–1130; Jones, Baylín: The epigenomics of cancer. *Cell*, 2007, 128:683–692; Qureshi, Mehler: Advances in epigenetics and epigenomics for neurodegenerative diseases. *Current neurology and neuroscience reports*, 2011, 11:464–473; Wood, Parsons, Jones, et al: The genomic landscapes of human breast and colorectal cancers. *Science*, 2007, 318:1108–1113.

Verkettbarkeit gerade bei Probanden mit Erkrankungsmarkern jedoch noch für sehr viel längere Zeiträume möglich sein.

3 Technische Gegenmaßnahmen

Als mögliche Gegenmaßnahmen für diese identifizierten Risiken für die Probanden von epigenetischen Studien möchten wir im Folgenden zwei verschiedene Verfahren vorschlagen. Die erste Gegenmaßnahme basiert darauf, nur eine Teilmenge der miRNAs zu veröffentlichen. In der Theorie können auf diese Weise die miRNAs unterdrückt oder versteckt werden, welche die Verknüpfungsangriffe ermöglichen. Die zweite Gegenmaßnahme basiert darauf, die Expressionsprofile zu verrauschen.

Da die Nutzbarkeit und Genauigkeit der Daten für die medizinische Forschung jedoch höchste Priorität haben, stehen die technischen Schutzmaßnahmen unter dem Vorbehalt, gleichzeitig die notwendige Nutzbarkeit zu erhalten, das heißt, Diagnosen mit Hilfe der Daten müssen noch möglich bleiben.

Unsere Evaluation der Gegenmaßnahmen basieren auf dem Worst-Case-Szenario, das wir im Rahmen unserer Angriffe identifiziert haben: Matching-Angriffe auf blutbasierte Proben (im Gegensatz zu plasmabasierten Proben). Zudem nehmen wir an, dass der Angreifer die bestmögliche Anzahl an Hauptkomponenten kennt.

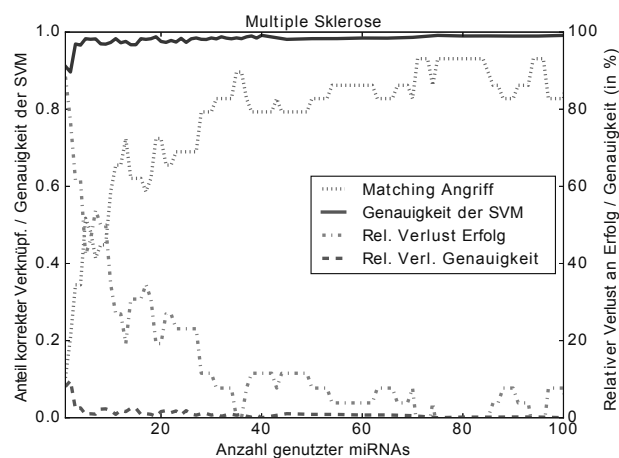
Den Nutzen der Daten ermitteln wir, indem wir mit dem in der Biomedizin üblichen Verfahren die Genauigkeit (*Accuracy*) einer Klassifizierung als gesund oder erkrankt bestimmen. Dazu wird üblicherweise eine radiale *Support Vector Machine* (SVM) genutzt. Eine geeignete Teilmenge an miRNAs wird durch das schrittweise Hinzufügen der miRNAs und erneutes Training der SVM ermittelt. Die Reihenfolge, in welcher die miRNAs als Features genutzt werden, wird meist mittels Signifikanzwerten bestimmt: Häufig sind dies *p*-Werte des *Wilcoxon-Mann-Whitney-Tests* (WMW). Die Genauigkeit ist definiert als die Anzahl korrekt klassifizierter Proben geteilt durch die Gesamtzahl an Proben. Wir nutzen den Mittelwert über eine zehnfache Cross-Validierung mit fünf Wiederholungen. Die *p*-Werte des WMW-Tests passen wir mit Hilfe der Benjamin-Hochberg-Anpassung an. Der relative Verlust des Nutzens der Daten wird dann als relative Veränderung der Genauigkeit in Bezug auf die ursprüngliche Genauigkeit gemessen, während der relative Verlust des Angriffserfolgs in Bezug auf den Maximalerfolg, d. h. 90%, angegeben wird.

3.1 Verstecken

Unsere erste Gegenmaßnahme ist das Verstecken einer Untergruppe der miRNAs. Dies hat den Vorteil, dass die Daten nicht perturbiert werden müssen, man also die korrekten Expressionswerte beibehalten kann. Wir evaluieren diese Gegenmaßnahme auf den blutbasierten Datensätzen, nachdem wir zuvor in einem Standardverfahren nicht aktive miRNAs gefiltert haben und nur solche miRNAs beibehalten, die in beiden Datensätzen vorhanden sind. Dies resultiert in 446 gemeinsamen miRNAs. Das Diagramm in Abb. 1 zeigt die Entwicklung des Angriffs und des Nutzens der Daten für einen Bereich von 1 bis 100 veröffentlichten miRNAs für Multiple Sklerose. Wir beschränken uns auf diesen Bereich, da die Genauigkeit der Klassifizierung in diesem Bereich bereits sehr hoch ist und die Erfolgsrate des Angriffs für mehr als

100 miRNAs bereits fast ihr Maximum erreicht. Wir veröffentlichen die miRNAs schrittweise in der Reihenfolge ihrer Signifikanzwerte. Dies erlaubt uns, vor allem die miRNAs mit geringeren Signifikanzwerten zu verstecken.

Abb. 1 | Verstecken von miRNA Expressionsprofilen am Beispiel Multiple Sklerose



Mit wenigen Ausnahmen liegt der relative Verlust an Genauigkeit immer unter 10%. Werden mehr als 20 miRNAs veröffentlicht, ist es zudem nicht möglich, den Angriffserfolg relativ gesehen um mehr als 50% zu verringern. Dennoch kann man häufig im Bereich zwischen drei und 20 miRNAs einen guten Trade-off finden. Anhand der Erkrankung Multiple Sklerose zeigen wir einen solchen Trade-off. Wir stellten fest, dass eine Menge von sieben miRNAs für Multiple Sklerose zu einer Verringerung des Erfolgs des Matching-Angriffs um 53,8% führt, während sich die Genauigkeit der Klassifizierung nur um 0,9% verschlechtert.

Dieses extensive Verstecken oder Unterdrücken von miRNA erscheint für bestimmte Zwecke durchaus sinnvoll.

Mit dem Verstecken von miRNA mit Ausnahme eines kleinen aber dafür hoch-akkuraten Datensets lassen sich im Hinblick auf den Nachweis konkreter bekannter Krankheiten gute Ergebnisse erzielen. Die Methode erscheint daher geeignet, um zum Beispiel miRNA-Daten zu Zwecken der Validierung von Forschungsergebnissen an andere Wissenschaftler veröffentlichen zu können, ohne die Probanden einem hohen Identifizierungsrisiko auszusetzen.

Allerdings ist es mit diesem Ansatz nicht mehr möglich, neue Assoziationen zwischen miRNAs und Erkrankungen zu finden, da nicht genug Daten zur Verfügung stehen. Für ergebnisoffene medizinische Forschung eignet sich diese Schutzmaßnahme daher nicht.

3.2 Verrauschen mit Differential Privacy

Unsere zweite Gegenmaßnahme basiert auf dem Verrauschen der Daten mithilfe von *Differential Privacy*. Bei dieser Gegenmaßnahme kann der Verantwortliche auf die miRNA-Expressionsprofile direkt zufälliges Rauschen addieren, bevor er zum Beispiel die gesamte Datenbank für Forschungszwecke zur Verfügung stellt. *Differential Privacy* eignet sich nicht nur für das Ver-

rauschen von Datenbankabfragen,¹⁴ sondern mit größerem Aufwand ebenso für das Verrauschen der gesamten Datenbank. Die Idee hinter dem Hinzufügen von Rauschen zu den rohen Expressionsdaten besteht darin, die Ununterscheidbarkeit zwischen verschiedenen Expressionsvektoren zu gewährleisten und somit die Verkettungs- und Identifizierungsmöglichkeiten des Angreifers zu reduzieren.

Der generalisierten Definition von Differential Privacy¹⁵ folgend, welche auch bereits erfolgreich für Standortdaten eingesetzt wurde,¹⁶ erreicht ein Mechanismus K *epigeno-indistinguishability* genau dann, wenn für alle m -miRNA-Expressionsprofile r_1 und r_2 gilt:

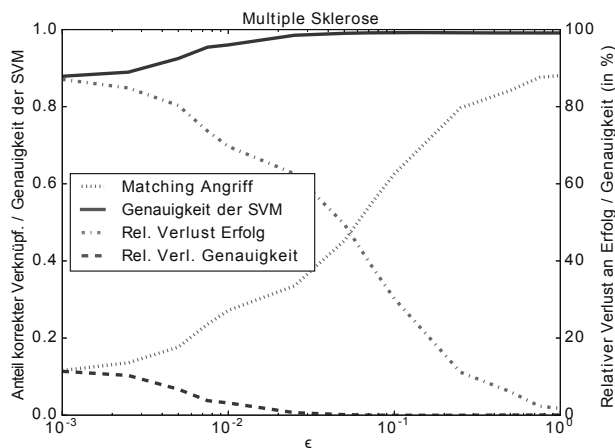
$$\Pr(K(r_1) \in \mathcal{S}) \leq \exp(\epsilon d_2(r_1, r_2)) \times \Pr(K(r_2) \in \mathcal{S})$$

wobei \mathcal{S} eine beliebige Teilmenge von möglichen Antworten und $d_2(\cdot, \cdot)$ die euklidische Distanz ist. Im Folgenden nehmen wir an, dass die Antworten im gleichen m -dimensionalen reellen Vektorraum \mathbb{R}^m liegen wie die ursprünglichen Profile. Wir erreichen *epigeno-indistinguishability*, indem wir die ursprünglichen Expressionsprofile durch Addieren des Vektors x verrauschen. x wird nach der Laplace-Wahrscheinlichkeitsdichtefunktion

$$g(x) = \frac{1}{\alpha} e^{-\epsilon \|x\|_1}$$

generiert, wobei α ein Normalisierungsfaktor ist, der garantiert, dass das Integral über g ergibt.

Abb. 2 | Verrauschen von miRNA Expressionsprofilen am Beispiel Multipler Sklerose



Mit Hilfe dieses Ansatzes verrauschen wir unsere Datensätze (basierend auf den gleichen 446 miRNAs wie zuvor) für verschiedene Werte für ϵ . Sowohl die Signifikanzwerte (p -Werte) als auch die SVM werden auf den verrauschten Daten ermittelt. Da das Verrauschen probabilistisch geschieht, wiederholen wir alle Experimente 50-mal und mitteln die Ergebnisse.

¹⁴ Leider betrachtet die Art. 29 Data Protection Working Party, Opinion 05/2014 on Anonymisation Techniques, 10. April 2014, S. 28, nur diesen sehr begrenzten Anwendungsfall.

¹⁵ Chatzikokolakis, Andrés, Bordenabe, Palamidessi: Broadening the scope of differential privacy using metrics. Privacy Enhancing Technologies, 2013, 82–102.

¹⁶ Andrés, Bordenabe, Chatzikokolakis, Palamidessi: Geo-indistinguishability: Differential privacy for location-based systems. Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security, 2013, 901–914.

Im Diagramm in Abb. 2 zeigen wir beispielhaft anhand der Erkrankung Multiple Sklerose, ähnlich wie beim Verstecken der miRNAs, den Trade-off zwischen Reduzierung der Identifizierbarkeit (Matchingangriff) und Genauigkeit der Daten für die Diagnose. Für die Multiple Sklerose bietet $\epsilon = 0,025$ den besten Trade-off mit einer Verschlechterung der Klassifizierungsgenauigkeit um nur 0,65% (Verstecken von miRNAs: 0,9%), während sich der Schutz vor Matching-Angriffen um 63% verbessert (Verstecken von miRNAs: 53,8%). Hier zeigt sich klar, dass insbesondere in diesem Fall das Verrauschen gegenüber Verstecken den besseren Trade-off bietet. Vergleicht man passende Trade-offs zwischen den beiden Gegenmaßnahmen Verstecken und Verrauschen, zeigt sich, dass bei 17 der 19 von uns untersuchten genetisch bedingten Erkrankungen das Verrauschen den besseren Datenschutz-Diagnosegenauigkeit-Trade-off bieten kann.¹⁷

3.3 Fazit

Es existieren *Privacy-Enhancing Technologies* (PETs), die auch bei großen Datenmengen eine signifikante Reduktion des Risikos einer Verkettung oder Identifizierung erreichen. Potentiell sind die vorgestellten Verfahren auch für eine Vielzahl weiterer Anwendungsfälle geeignet.

Leider gehen die vorgeschlagenen Maßnahmen jedoch jeweils mit einem Trade-off zwischen Datenschutz und Genauigkeit und Nutzbarkeit der epigenetischen Daten einher. Wo je nach Einzelfall dieser Trade-off zu ziehen ist, muss sich nach dem Zweck der Verarbeitung, der Rechtsgrundlage und den Risiken für die Betroffenen bemessen.

Auch wenn unsere Ergebnisse immer noch eine erstaunliche Genauigkeit der Diagnose erlauben, ist es unklar wann gerade im Bereich der medizinischen Forschung eine solche Ungenauigkeit tolerierbar ist und sein muss. In Anbetracht dessen, dass die medizinische Forschung bestrebt ist, so exakte Diagnose-Ergebnisse wie möglich zu erzielen, wird sich zeigen müssen, ob sie davon zu überzeugen ist, zum Schutz der Betroffenen in einigen Forschungsvorhaben einen nicht zu vernachlässigenden Nutzenverlust in Kauf zu nehmen.

4 Rechtliche Bewertung

Aus datenschutzrechtlicher Perspektive stellt sich für die vorgestellten Maßnahmen die Frage, wie diese einzuordnen sind.

Im Hinblick auf die Verarbeitung personenbezogener Daten zu Forschungszwecken macht Art. 89 Abs. 1 DSGVO Vorgaben zu geeigneten Garantien bei der Verarbeitung, die die Rechte und Freiheiten der Betroffenen in geeigneter Weise schützen. Die DSGVO verweist dabei insbesondere auf Maßnahmen der Datenminimierung (Satz 2) durch Pseudonymisierung (Satz 3) oder vorrangig durch Anonymisierung (Satz 4). Dem Verantwortlichen obliegt es demzufolge bereits während der Planung eines Forschungsprojekts zu prüfen, ob die Verarbeitungszwecke auch mit anonymisierten oder wenigstens pseudonymisierten Daten erfüllt werden können.¹⁸ Während die Anonymisierung Daten vom Anwendungsbereich der DSGVO ausnimmt, stellt

¹⁷ Details dieser Analysen finden sich in unserer Veröffentlichung Backes, Berrang, Hecksteden, Humbert, Keller, Meyer: Privacy in epigenetics: Temporal linkability of microRNA expression profiles. USENIX Security 2016.

¹⁸ Raum, in Ehmann/Selmayr Datenschutz-Grundverordnung, Art. 89 Rn. 35, 37.

eine Pseudonymisierung lediglich eine Maßnahme der Datenminimierung dar.

In jedem Fall handelt es sich bei den vorgeschlagenen Verfahren um Maßnahmen der Datenminimierung. Die interessantere rechtswissenschaftliche Fragestellung ist jedoch, ob die Verfahren eventuell den hohen Anforderungen, die die EU-Datenschutzgrundverordnung (DSGVO) an eine Anonymisierung stellt, genügen können.

4.1 Grundsätzliche Möglichkeit der Anonymisierung

Die DSGVO geht nicht auf die Eignung konkreter Anonymisierungsverfahren ein. Vielmehr ist Art. 4 Nr. 1 DSGVO in Verbindung mit Erwägungsgrund 26 maßgeblich dafür, ob eine „echte Anonymisierung“¹⁹ der Daten vorliegt. Dies ist der Fall, wenn keine direkte oder indirekte (Erwägungsgrund 26 nennt hier „Aussondern“ als Beispiel) Identifizierbarkeit der Betroffenen mehr möglich ist. Obwohl der Erwägungsgrund 26 auf eine risikobasierte Einschätzung einer möglichen Re-Identifizierung abzielt – „nach allgemeinem Ermessen wahrscheinlich“ – geht die herrschende Literaturmeinung von einem binären Begriff der Anonymität aus. Daten können nur entweder anonym oder personenbeziehbar sein.²⁰

Die Art. 29 Arbeitsgruppe konkretisierte die Anforderungen an Anonymisierung im Jahr 2014 noch nach Maßgabe der EU-Datenschutzrichtlinie. Eine effektive Anonymisierung verhindere, dass alle Beteiligten eine Person in einem Datensatz ausfindig machen, zwei Datensätze innerhalb eines Datensatzes (oder zwischen zwei getrennten Datensätzen) miteinander verknüpfen und daraus Informationen ableiten (inferieren) können.²¹ „Alle Beteiligten“ schließe nach Auffassung der Datenschützer explizit auch den Verantwortlichen selbst mit ein; solange dieser oder irgendein Dritter ein identifizierbares Datenset behält, sei eine Anonymisierung eines daraus entstandenen Datensets rechtlich nicht möglich.²² Diese Interpretation wird von einigen Kommentatoren in gleicher Weise für die DSGVO fortgesetzt.²³

In Anbetracht der langen gesetzlichen Speicherfristen für medizinische und insbesondere genetische (Roh-)Daten in den EU-Mitgliedsstaaten, würde dies zu dem Ergebnis führen, dass eine Anonymisierung dieser Daten für Forschungszwecke während der Speicherfrist gar nicht möglich sei. Da es sich bei den Betroffenen häufig um Patienten in einer andauernden Behandlung handelt, verlängert sich dieser Zeitraum sogar auf die Dauer der Behandlung plus die gesetzliche Speicherfrist.

Die Option einer zumindest rechtlichen Unmöglichkeit der Anonymisierung wird in der medizinischen Forschung noch kaum diskutiert. Auch wenn eine zwingende Anwendbarkeit der DSGVO auf jegliche medizinischen Forschungsdaten auf den ersten Blick zum Schutz der Betroffenen vorteilhaft erscheinen mag, gingen aus unserer Sicht wichtige Anreize verloren, die Personenbeziehbarkeit und damit einen Großteil der Risiken durch Technikgestaltung tatsächlich so weit wie möglich zu reduzieren. Viel-

mehr wäre zu erwarten, dass Verantwortliche und Forscher sich auf ein „gut genug“ an Pseudonymisierung ergänzt mit pauschalen Einwilligungen der Probanden beschränken würden.

Aus unserer Sicht ist diese Interpretation der DSGVO jedoch keineswegs zwingend. Der rechtswissenschaftliche Diskurs über einen absoluten vs. einen relativen Begriff des Personenbezugs erscheint auch durch den Wortlaut des Erwägungsgrund 26 nicht abschließend geklärt.²⁴ Der EuGH folgte in der Breyer Entscheidung 2016 einem Mittelweg. Es sei nicht jedes Wissen eines realen oder hypothetischen Dritten relevant, zu prüfen sei jedoch, ob derjenige im Besitz der Daten Mittel und Möglichkeit habe, die „vernünftigerweise zur Bestimmung der betreffenden Person eingesetzt werden“ können.²⁵ Der EuGH bejahte dies für den Webseiten-Anbieter, der über rechtliche Mittel verfüge, um über identifizierbare Betroffene hinter dynamischen IP-Adressen Auskunft beim Provider zu verlangen.

Folgt man diesem Mittelweg des EuGH, der auch unter den Vorgaben der DSGVO Bestand behalten dürfte, bleibt eine Anonymisierung von medizinischen Forschungsdaten trotz einer Pflicht zur Aufbewahrung der Rohdaten möglich. Es käme entscheidend darauf an, ob die Stelle, die das anonymisierte Datenset erhält „vernünftigerweise“ und „nach allgemeinem Ermessen wahrscheinlich“ Mittel für eine Re-Identifizierung besitzt. Bei dieser Prüfung sind die in Erwägungsgrund 26 DSGVO genannten Kriterien, insbesondere die zum Zeitpunkt der Verarbeitung verfügbare Technologie und technologische Entwicklung zu berücksichtigen. Diese Prüfung muss im Einzelfall erfolgen. Eine effektive Anonymisierung genetischer Forschungsdaten zum Zweck der Weitergabe an andere Wissenschaftler oder Veröffentlichung erscheint aber zumindest nicht grundsätzlich rechtlich ausgeschlossen.

4.2 Eignung der vorgeschlagenen Maßnahmen

Im Hinblick auf die vorgestellten Verfahren erweist sich die rechtliche Dichotomie von personenbeziehbar und anonym als problematisch. Aus Sicht der Informatik ist die Reduzierung der Identifizierbarkeit und damit die Verringerung des Risikos einer Re-Identifizierbarkeit ein dynamischer, optimierbarer Prozess. Dies erkennt auch die Art. 29 Arbeitsgruppe an.²⁶

Bei den von uns vorgeschlagenen Maßnahmen stellt sich in der praktischen Anwendung die Frage, wann die Möglichkeit der Identifizierbarkeit, Verkettbarkeit oder Inferenz im Sinne der Risiko-Abschätzung hinreichend gering sind. Die Konfidenz mit der ein Angreifer eine bestimmte Aussage treffen oder eine Re-Identifizierung vornehmen kann, lässt sich mathematisch bestimmen. Sie wird aber nicht (oder nur in den seltensten Fällen) 0 betragen.

Anschaulich wird dies am folgenden Beispiel: Angreifer A besitzt das Zusatzwissen, dass sein Nachbar N an einem bestimmten Datum in die Notaufnahme eingeliefert wurde. Das Krankenhaus publiziert statistische Daten und weist für diesen Tag drei aufgenommene Patienten mit den Attributen Myokardinfarkt, Pneumonie und Schlaganfall aus. A kann somit mit einer Konfidenz von 33,3% eine Aussage über Ns akute Erkrankung machen.

19 *Klabunde*, in Ehmann/Selmayr Datenschutz-Grundverordnung, Art. 4 Rn. 16.

20 *Karg*, DuD 2015, 520 (523); *Klar/Kühling*, in Kühling/Buchner Datenschutz-Grundverordnung, Art. 4 Nr. 1 Rn 32.; *Klabunde*, in Ehmann/Selmayr Datenschutz-Grundverordnung, Art. 4 Rn. 16.

21 *Art. 29 Data Protection Working Party*, Opinion 05/2014 on Anonymisation Techniques, 10 April 2014, S. 9.

22 *Art. 29 Data Protection Working Party*, Opinion 05/2014 on Anonymisation Techniques, 10 April 2014, S. 9.

23 *Klabunde*, in Ehmann/Selmayr Datenschutz-Grundverordnung, Art. 4 Rn. 16.

24 *Marnau*, DuD 2016, 428 (429 f.).

25 EuGH, 19.10.2016 – C-582/14, Rn 45.

26 *Art. 29 Data Protection Working Party*, Opinion 05/2014 on Anonymisation Techniques, 10. April 2014, S. 4.

Die Anwendung von Differential Privacy führt zu ähnlichen Fragestellungen. Das Verrauschen der Daten um einen zuvor ermittelten Faktor erlaubt dem Angreifer immer noch zumindest den Korridor zu erkennen, in dem sich der echte Wert befindet.

Mit der letztlichen Entscheidung, z. B. im Rahmen einer Datenschutzfolgenabschätzung, ob eine Maßnahme ausreichend ist, um das Risiko einer Identifizierbarkeit, Verkettbarkeit oder Aussonderung einzelner Betroffener nach allgemeinem Ermessen auszuschließen, werden Verantwortliche derzeit jedoch allein gelassen. Weder Gerichte, noch Aufsichtsbehörden oder die Forschung haben hierfür bereits in ausreichendem Maße Kriterien oder Best Practices entwickelt.

Diesem Fehlen an Kriterien versuchen Firmen, die innovative Anonymisierungskonzepte einsetzen, jetzt auf anderer Weise zu begegnen. Ähnlich wie es im Bereich der IT-Sicherheit bereits Penetrationstests und *Bug-Bounty*-Programme als etablierte Verfahren gibt, stellen sie sich einem Re-Identifizierungs-Testen mit entsprechendem Bounty-Programm.²⁷ Selbstverständlich ist bei solchen Programmen bei der Ausgestaltung darauf zu achten, die Rechte und Freiheiten der Betroffenen auch vor gutmeinenden Testern zu schützen, zum Beispiel im Rahmen eines geschlossenen Kreises von Teilnehmern und vertraglichen Bedingungen für Teilnahme. Dennoch sind solche Programme aus unserer Sicht ausdrücklich zu begrüßen, um es zu ermöglichen langfristig belastbare Kriterien für die Bewertung nach Erwägungsgrund 26 DSGVO zu identifizieren. Denkbar wäre es, ein solches Testprogramm auch als Maßnahme in die Datenschutzfolgenabschätzung zu integrieren und regelmäßig zu wiederholen, um die technische Entwicklung abzubilden.

In diesem Kontext betrachten wir mit Sorge die vermehrten Vorstöße von Politikern und Aufsichtsbehörden, Versuche der Re-Identifizierung vermeintlich anonymisierter Daten zu untersagen oder sogar unter Strafe zu stellen. Zum Thema Big Data und Gesundheitsdaten findet sich zum Beispiel in den Forderungen der Konferenz der unabhängigen Datenschutzbehörden des Bundes und der Länder für die neue Legislaturperiode: „Verknüpfungen zwischen verschiedenen Datenbeständen, die Gesundheitsdaten enthalten, dürfen nur auf der Grundlage spezieller rechtlicher Regelungen zugelassen werden. Die Re-Identifizierung und

unerlaubte Zusammenführung von Daten [...] sind zu verbieten und unter Strafe zu stellen.“²⁸

Ein solcher Straftatbestand würde den Bereich der Forschung und Entwicklung von Anonymisierungsverfahren in Deutschland inkriminieren, ohne auf internationaler Ebene einen besseren Schutz für die Betroffenen zu erreichen. Vielmehr würde eine solche Kriminalisierung von Zusammenführungs- und Re-Identifizierungsversuchen sogar zu einem Weniger an Schutz führen, da neue Re-Identifizierungsrisiken, wie wir sie für epigenetische Daten identifiziert haben, ggf. nicht (rechtzeitig) entdeckt und Technologien nicht mehr weiterentwickelt werden.

5 Fazit

Die Vereinbarkeit medizinischer Forschung zum gesamtgesellschaftlichen Nutzen mit den Rechten und Freiheiten der betroffenen Probanden wird auch in den kommenden Jahren eine Kernherausforderung für Mediziner, Aufsichtsbehörden und Wissenschaftler vieler Disziplinen bleiben. Auch wenn vielversprechende technische Verfahren für einen besseren Datenschutz existieren, sind noch zahlreiche Fragen nach dem Ausgleich zwischen Datenschutz und Datennutzbarkeit unbeantwortet. Während die Datenschutzgrundverordnung in Zukunft höhere technisch-organisatorische Anforderungen an Datenverarbeitung zu Forschungszwecken stellen wird (und somit begrüßenswerter Weise vom Vorrang der Einwilligung abrückt), führt das Fehlen konkreter Maßstäbe und Richtlinien zu einer nur schwachen Incentivierung für den Einsatz von Anonymisierungsverfahren in der Praxis.

Literatur

- [1] *Backes, Berrang, Hecksteden, Humbert, Keller, Meyer: Privacy in epigenetics: Temporal linkability of microRNA expression profiles. USENIX Security 2016.*
- [2] *Marnau: Anonymisierung, Pseudonymisierung und Transparenz für Big Data, DuD 2016, 428.*
- [3] *Herbst: Rechtliche und ethische Probleme des Umgangs mit Proben und Daten bei großen Biobanken, DuD 2016, 371.*

²⁸ Konferenz der unabhängigen Datenschutzbehörden des Bundes und der Länder: Grundsatzpositionen und Forderungen für die neue Legislaturperiode, https://datenschutz.saarland.de/uploads/media/DSK_Grundsatzpositionen-und-Forderungen-fuer-die-neue-Legislaturperiode.pdf.

²⁷ Aircloak Challenge <https://aircloak.com/challenge/>.