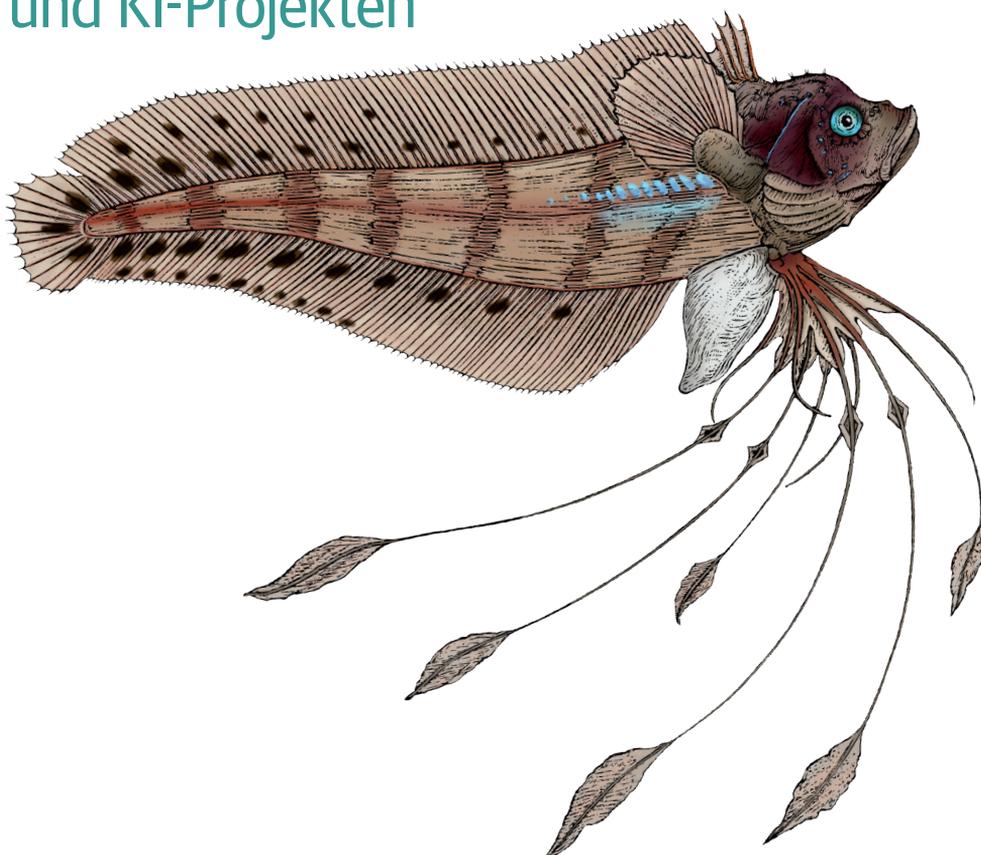


O'REILLY®

Inklusive
DSGVO

Data Privacy in der Praxis

Datenschutz und Sicherheit in Daten-
und KI-Projekten



Katharine Jarmul

Übersetzung von Marcus Fraaß

Stimmen zum Buch

Data Privacy in der Praxis

Ein absolutes Standardwerk für alle Flughöhen, vom Data Scientist bis zur Verbraucherschützerin. Katharine Jarmul erklärt verständlich und umfassend die Theorie und stellt praxisnahe Beispiele und Lösungsansätze vor. Ein pointiertes, leidenschaftliches Plädoyer für ein im Alltag und in Unternehmen viel zu oft vernachlässigtes Thema: Data Privacy.

– Karma Lüdtko, IT Innovations bei Bundesdruckerei und Expert:in
Cloud-Infrastrukturen

Katharine Jarmuls *Data Privacy in der Praxis* ist ein ausgezeichnete praktischer Leitfaden für die Integration des Datenschutzes in digitale Produkte. Ein Muss für alle, die digitale Produkte entwickeln.

– Alexander CS Hendorf, Gründer opotoc GmbH,
Vorsitzender Python Softwareverband e.V.

Katharine Jarmul bietet mit diesem Buch eine hilfreiche und praxisnahe Handreichung für Data Scientists, Privacy Engineers und Softwareentwickler*innen, die sich mit den technischen, juristischen und ethischen Grundlagen von Datensicherheit und Security auseinandersetzen wollen. Die vielen Beispiele und hilfreichen Erklärungen ermöglichen es Entwickler*innen, ihr neu gewonnenes Wissen direkt in ihre Arbeit zu integrieren. Damit schafft die Autorin Zugang zu zukunftsweisenden Fähigkeiten, die immer weiter an Relevanz gewinnen werden – für diejenigen, die an den Systemen von morgen arbeiten. Ein wichtiger Beitrag für eine wünschenswerte und gerechte Digitalisierung, in der »Privacy by Design« der Standard ist.

– Fiona Krakenbürger, Co-Founder und
Chief Technology Officer Sovereign Tech Fund

Es war noch nie so leicht, KI-Prototypen zu entwickeln, aber die Lösung von realen Problemen mit realen Daten ist nach wie vor eine große Herausforderung. Katharine Jarmul bietet in *Data Privacy in der Praxis* einen umfassenden und anschaulichen Überblick über Datensicherheit und Privatsphäre für KI und Machine Learning, mit technisch fundierten Beispielen und verschiedenen Lösungsansätzen aus der Praxis. Die vorgestellten und direkt anwendbaren Methoden ermöglichen es jedem, selbstbestimmt mit privaten Daten zu arbeiten und transparente und zukunftsfähige KI-Systeme zu entwickeln.

– Ines Montani, Gründerin und CEO von Explosion und Entwicklerin von spaCy

Data Privacy geht uns alle etwas an. Gerade in Zeiten, wo Generative AI und LLMs auch im Mainstream angekommen sind, sollten alle die Wichtigkeit von Privatsphäre und vom verantwortungsvollen Handeln mit persönlichen Daten erkannt haben. Katharine zeigt uns in diesem Buch nicht nur die Risiken und Konsequenzen von allem, was mit Data Privacy schief laufen kann, sondern auch, wie wir anhand von praxisorientierten Beispielen aus der Realität diesen Risiken vorbeugen können.

Es ist nicht nur ein wichtiges Buch, das alle Data Professionals als Grundlektüre lesen sollten – es macht auch dank Katharines zugänglichem Schreibstil und Ihrer unbezweifelbaren Leidenschaft für dieses Thema unheimlich Spaß zu lesen!

– Tiankai Feng
Data Strategy & Data Governance Lead, Thoughtworks Europe

Will man ein Softwareprodukt mit künstlicher Intelligenz bauen, so steht die Datenbeschaffung ganz am Anfang des Projektes, und bereits in dieser Projektphase liegt die vermeintlich größte Hürde. Daten sind aufgrund von Datenschutzbedenken schwer zu beschaffen. Dabei gibt es unzählige Möglichkeiten, KI-Software datenschutzkonform zu bauen. Die Kombination von klassischer Informationssicherheit mit dem aktuellen Stand der Wissenschaft ermöglicht sogar einen Privacy-by-Design-Ansatz. Katharine Jarmul stellt in ihrem Buch *Data Privacy in der Praxis* diesen Ansatz vor und gibt einen umfassenden Überblick über die Details der damit verbundenen Themen. Dies ermöglicht es, Datenschutz nicht als Hindernis, sondern als technologische Basis zu verstehen, die regulatorische Aspekte berücksichtigt und die Privatsphäre jedes Einzelnen respektiert. Das Buch deckt die verschiedenen Aspekte der Daten- und Informationssicherheit in Softwareprojekten ab und eignet sich daher hervorragend für alle Mitarbeiter:innen in Softwareprojekten: Entwickler:innen, Projektleiter:innen und auch Sicherheitsexpert:innen.

– Dr. Maria Börner
*Competence Center AI Lead bei Westernacher Solutions,
Partnership Germany Lead bei Women in AI*

Data Privacy in der Praxis bietet genau das, was der Titel verspricht – eine praxisorientierte Darstellung der verschiedenen Ansätze im Bereich des Datenschutzes, die auch den wirtschaftlichen Nutzen im Zusammenhang mit der Nutzung personenbezogener Daten ausreichend berücksichtigt.

– Rebecca Parsons, *Chief Technology Officer bei Thoughtworks*

Die Datenlandschaft wird mit jedem Jahr komplexer. Der Druck der Regulierungsbehörden in Bezug auf Datenschutz und Datensouveränität sowie Transparenz, Erklärbarkeit und Fairness von Algorithmen nimmt weltweit zu. Es ist schwieriger denn je, Daten intelligent zu verwalten. Aber die Werkzeuge zur Bewältigung dieser Herausforderungen sind besser denn je, und dieses Buch ist eines dieser Werkzeuge. Katharines praxisorientierte, pragmatische und umfassende Behandlung des Themas *Data Privacy* ist genau das richtige Buch, um die Herausforderungen der 2020er-Jahre und darüber hinaus zu meistern. Sie schafft es, ihre fachlich fundierten Ausführungen mit leicht verständlichen Übersichten über die neuesten technologischen Ansätze und Architekturen zu verbinden. Dieses Buch ist für jeden nützlich, vom CDO bis zum Data Analyst und jedem dazwischen

– Emily F. Gorcenski, *Principal Data Scientist und Data & AI Service Line Lead bei Thoughtworks*

Manche Data Scientists sehen Datenschutz als etwas an, das sie in ihrer Arbeit eintrübt. Wenn Sie jedoch nicht zu dieser Gruppe gehören, wenn Sie glauben, dass Datenschutz sowohl in moralischer als auch in wirtschaftlicher Hinsicht erstrebenswert ist, wenn Sie die Stringenz und die Möglichkeiten des Privacy Engineering schätzen, wenn Sie sich über den aktuellen Stand der Technik auf diesem Gebiet informieren möchten, dann ist dieses Buch genau das Richtige für Sie.

– Chris Ford, *Head of Technology bei Thoughtworks Spanien*

Copyright und Urheberrechte:

Die durch die dpunkt.verlag GmbH vertriebenen digitalen Inhalte sind urheberrechtlich geschützt. Der Nutzer verpflichtet sich, die Urheberrechte anzuerkennen und einzuhalten. Es werden keine Urheber-, Nutzungs- und sonstigen Schutzrechte an den Inhalten auf den Nutzer übertragen. Der Nutzer ist nur berechtigt, den abgerufenen Inhalt zu eigenen Zwecken zu nutzen. Er ist nicht berechtigt, den Inhalt im Internet, in Intranets, in Extranets oder sonst wie Dritten zur Verwertung zur Verfügung zu stellen. Eine öffentliche Wiedergabe oder sonstige Weiterveröffentlichung und eine gewerbliche Vervielfältigung der Inhalte wird ausdrücklich ausgeschlossen. Der Nutzer darf Urheberrechtsvermerke, Markenzeichen und andere Rechtsvorbehalte im abgerufenen Inhalt nicht entfernen.

Data Privacy in der Praxis

*Datenschutz und Sicherheit in
Daten- und KI-Projekten*

Katharine Jarmul

*Deutsche Übersetzung von
Marcus Fraaß*

O'REILLY®

Katharine Jarmul

Lektorat: Alexandra Follenius

Übersetzung: Marcus Fraaß

Korrektorat: Sibylle Feldmann, www.richtiger-text.de

Satz: III-satz, www.drei-satz.de

Herstellung: Stefanie Weidner

Umschlaggestaltung: Karen Montgomery, Michael Oréal, www.oreal.de

Bibliografische Information der Deutschen Nationalbibliothek

Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über <http://dnb.d-nb.de> abrufbar.

ISBN:

Print 978-3-96009-233-9

PDF 978-3-96010-816-0

ePub 978-3-96010-817-7

1. Auflage 2024

Translation Copyright für die deutschsprachige Ausgabe © 2024 dpunkt.verlag GmbH

Wiebinger Weg 17

69123 Heidelberg

Authorized German translation of the English edition of *Practical Data Privacy* ISBN 9781098129460

© 2023 Kjamistan, Inc. This translation is published and sold by permission of O'Reilly Media, Inc., which owns or controls all rights to publish and sell the same.

Dieses Buch erscheint in Kooperation mit O'Reilly Media, Inc. unter dem Imprint »O'REILLY«.

O'REILLY ist ein Markenzeichen und eine eingetragene Marke von O'Reilly Media, Inc. und wird mit Einwilligung des Eigentümers verwendet.

Schreiben Sie uns:

Falls Sie Anregungen, Wünsche und Kommentare haben, lassen Sie es uns wissen: komentar@oreilly.de.

Die vorliegende Publikation ist urheberrechtlich geschützt. Alle Rechte vorbehalten. Die Verwendung der Texte und Abbildungen, auch auszugsweise, ist ohne die schriftliche Zustimmung des Verlags urheberrechtswidrig und daher strafbar. Dies gilt insbesondere für die Vervielfältigung, Übersetzung oder die Verwendung in elektronischen Systemen.

Es wird darauf hingewiesen, dass die im Buch verwendeten Soft- und Hardware-Bezeichnungen sowie Markennamen und Produktbezeichnungen der jeweiligen Firmen im Allgemeinen warenzeichen-, marken- oder patentrechtlichem Schutz unterliegen.

Alle Angaben und Programme in diesem Buch wurden mit größter Sorgfalt kontrolliert. Weder Autorin noch Übersetzer noch Verlag können jedoch für Schäden haftbar gemacht werden, die in Zusammenhang mit der Verwendung dieses Buchs stehen.

Vorwort	15
Einleitung	19
1 Data Governance und einfache Datenschutzansätze	35
Data Governance: Was ist das?	36
Sensible Daten identifizieren	39
Persönlich identifizierende Informationen (PII) identifizieren	42
Datennutzung dokumentieren	43
Grundlagen der Datendokumentation	44
Unbekannte Daten aufspüren und dokumentieren	49
Data-Lineage-Tracking	52
Versionskontrolle für Daten	55
Grundlegender Datenschutz: Pseudonymisierung beim Privacy by Design	58
Zusammenfassung	63
2 Anonymisierung	65
Was ist Anonymisierung?	65
Definition von Differential Privacy	68
Das Epsilon verstehen: Was ist der Privacy Loss?	70
Was Differential Privacy garantiert und was nicht	73
Differential Privacy verstehen	75
Differential Privacy in der Praxis: Anonymisierung der Zensusdaten in den USA	75
Differential Privacy auf Basis des Laplace-Mechanismus	78
Differential Privacy auf Basis des Laplace-Mechanismus: ein simpler Ansatz	81
Sensitivität und Fehler	83
Privacy Budgets und deren Aufteilung	85

Weitere Mechanismen erkunden: Differential Privacy mittels des gaußschen Rauschens	88
Laplace-verteilt und gaußsches Rauschen im Vergleich	90
Differential Privacy in der Praxis: Debiasing von Differential- Privacy-Ergebnissen	94
Sensitivität und Privacy Units	95
Wie steht es mit k-Anonymity?	96
Zusammenfassung	99
3 Datenschutz in Datenpipelines integrieren	101
Datenschutz in Datenpipelines integrieren	101
Geeignete Datenschutzmaßnahmen konzipieren	102
Die Nutzerinnen und Nutzer besser einschätzen können	104
Datenschutz in Datenpipelines integrieren	105
Testen und validieren	106
Datenschutz und Data Governance in Pipelines integrieren	107
Ein Beispiel für einen Workflow zur gemeinsamen Nutzung von Daten	107
Informationen zur Datenherkunft und Einwilligung im Rahmen der Datenerhebung zusätzlich erfassen	110
Differential-Privacy-Bibliotheken in Pipelines verwenden	114
Daten anonymisiert erheben	119
Datenerhebung unter Anwendung von Differential Privacy bei Apple	119
Warum bei Chrome der ursprüngliche Differential-Privacy-Ansatz im Rahmen der Datenerhebung eingestellt wurde	122
Zusammenarbeit mit dem Data-Engineering-Team und Führungskräften	125
Verantwortung teilen	126
Workflows zur Dokumentation von Datenschutzmaßnahmen und -empfehlungen erstellen	127
Datenschutz als zentrales Wertversprechen	127
Zusammenfassung	128
4 Angriffe auf die Privatsphäre	131
Angriffe auf die Privatsphäre: eine Analyse gängiger Angriffsvektoren . . .	131
Der Netflix-Prize-Angriff	131
Linkage Attacks	134
Singling Out Attacks	137
Der Strava-Heat-Map-Angriff	138
Membership Inference Attack	141
Auf sensible Merkmale zurückschließen	144
Andere Leakage Attacks auf Modelle: Memorierung	146

Data Exfiltration Attacks auf ChatGPT und andere LLMs	147
Model-Stealing Attacks	150
Informationen aus Prompts und zusätzlichen Dokumenten extrahieren	152
Angriffe auf Privacy-Mechanismen	153
Datensicherheit	155
Zugriffskontrolle	157
Schutz vor Datenverlust	157
Zusätzliche Sicherheitsvorkehrungen	158
Threat Modeling und Incident-Response-Pläne	159
Angriffe mithilfe von Eintrittswahrscheinlichkeiten bewerten	160
Ein »durchschnittlicher« Angreifer	160
Risiken bewerten und Bedrohungen einschätzen	162
Vorkehrungen für die Datensicherheit, die auch dem Schutz der Privatsphäre dienen können	163
Die Websicherheit-Basics anwenden	164
Trainingsdaten und Modelle schützen	164
Über neue Angriffe auf dem Laufenden bleiben	166
Zusammenfassung	167
5 Machine Learning und Data Science datenschutzkonform gestalten	169
Privacy-preserving Machine Learning (PPML)	170
Techniken zur Wahrung der Privatsphäre in einem typischen Data-Science- bzw. ML-Workflow	170
Privacy-preserving Machine Learning in der Praxis	175
Stochastisches Gradientenabstiegsverfahren mit Differential Privacy (DP-SGD)	176
Open-Source-Bibliotheken für PPML	179
Differential Privacy bei LLMs und vergleichbaren generativen Systemen anwenden	183
Feature Engineering mit Differential Privacy	185
Einfachere Methoden anwenden	188
Machine Learning dokumentieren	189
Andere Wege, um die Privatsphäre beim Machine Learning zu schützen	192
Datenschutz in die Architektur für Daten- und Machine-Learning- Projekte integrieren	196
Ihre Datenschutzerfordernungen verstehen	196
Monitoring des Datenschutzes	198
Zusammenfassung	200

6	Federated Learning und Data Science	201
	Verteilte Daten	201
	Warum verteilte Daten nutzen?	202
	Wie funktioniert die verteilte Datenanalyse?	204
	Datenschutz bei verteilten Daten mittels Differential Privacy gewährleisten	207
	Federated Learning	209
	Die Entwicklung des Federated Learning im Überblick	210
	Weshalb, wann und wie Sie Federated Learning einsetzen sollten ...	212
	Federated-Learning-Systeme konzipieren	215
	Mögliche Arten des Deployments	216
	Potenzielle Sicherheitsrisiken	219
	Anwendungsbereiche	221
	Deployment mit Federated-Learning-Bibliotheken und -Tools	222
	Open-Source-Bibliotheken für Federated Learning	223
	Flower: eine Federated-Learning-Bibliothek für verschiedene Open-Source-Backends	224
	Federated Data Science – ein Ausblick	227
	Zusammenfassung	228
7	Encrypted Computation	229
	Was genau ist Encrypted Computation?	229
	Wann Encrypted Computation verwendet werden sollte	230
	Unterschied zwischen Datenschutz und Geheimhaltung	232
	Threat Modeling	234
	Verschiedene Arten der Encrypted Computation	236
	Secure Multiparty Computation	236
	Homomorphe Verschlüsselung	247
	Reale Anwendungsfälle im Zusammenhang mit Encrypted Computation	256
	Private Set Intersection	256
	Private Join and Compute	259
	Sichere Aggregation (Secure Aggregation)	260
	Encrypted Machine Learning	262
	Die ersten Schritte mit PSI und Moose	264
	Vision einer Welt mit sicherem Datenaustausch	272
	Zusammenfassung	273
8	Datenschutzrechtliche Aspekte	275
	Die DSGVO im Überblick	276
	Grundlegende Rechte nach DSGVO	276
	Datenverantwortlicher und Datenverarbeiter – eine Abgrenzung ...	279
	Technologien zur Verbesserung des Datenschutzes (PETs) im Hinblick auf die DSGVO einsetzen	281

Die Datenschutz-Folgenabschätzung der DSGVO: agile und iterative Risikobewertung	284
Recht auf Erläuterung: Nachvollziehbarkeit und Datenschutz	289
Der California Consumer Privacy Act (CCPA)	289
Technologien zur Verbesserung des Datenschutzes (PETs) im Hinblick auf den CCPA einsetzen	291
Weitere Vorschriften: HIPAA, LGPD, PIPL und andere	292
Datenschutzrechtliche Aspekte des AI Act	294
Data Governance Act	296
Data Act	297
Interne Richtlinien und Verträge	298
Datenschutzrichtlinien und Nutzungsbedingungen lesen	298
Auftragsverarbeitungsverträge lesen	301
Richtlinien, Leitfäden und Verträge lesen	302
Zusammenarbeit mit Rechtsexperten	303
Einhaltung von vertraglichen Vereinbarungen und Vertragsrecht	304
Datenschutzbestimmungen auslegen	305
Unterstützung und Rat einholen	306
Gemeinsam Definitionen und Ideen erarbeiten	307
Technische Beratung leisten	307
Data Governance 2.0	308
Was ist Federated Governance?	309
Eine Kultur des Experimentierens fördern	311
Den Schutz der Privatsphäre (PETs) verbessern mit funktionierender Dokumentation und Plattformen mit integrierten Technologien	312
Zusammenfassung	313
9 Datenschutz und Anwendungen aus der Praxis	315
Datenschutz- und Sicherheitsrisiken in der Praxis managen	316
Datenschutzrisiken bewerten und managen	316
Mit Ungewissheit umgehen und gleichzeitig für die Zukunft planen	319
Der Einsatz von Datenschutztechnologien in der Praxis: eine Analyse konkreter Anwendungsfälle	322
Federated Marketing: Marketingkampagnen unter Berücksichtigung des Datenschutzes durchführen	322
Public-Private-Partnerships: gemeinsame Nutzung von Daten im öffentlichen Gesundheitsdienst	326
Machine Learning mit anonymisierten Daten: DSGVO-konforme Lösungen in einem iterativen Trainings-Setting	329
Business-to-Business-Anwendung: Zugriff auf Daten aus erster Hand	331

Schrittweise Integration und Automatisierung von Datenschutz im Rahmen von Machine Learning	333
Iterative Erkundung	334
Datenschutzanforderungen dokumentieren	335
Ansätze evaluieren und kombinieren	338
Prozesse zunehmend automatisieren	340
Datenschutz zur Normalität werden lassen	341
Den Weg in die Zukunft ebnen: mit Forschungsbibliotheken arbeiten und Forschungsgruppen einbeziehen	342
Mit externen Forscherinnen und Forschern zusammenarbeiten	343
In interne Forschung investieren	344
Zusammenfassung	346
10 Häufig gestellte Fragen und ihre Antworten!	347
Encrypted Computation und Confidential Computing	347
Ist Secure Computation quantensicher?	348
Kann ich Enklaven verwenden, um Datenschutzprobleme oder Probleme im Zusammenhang mit der Geheimhaltung von Daten zu lösen?	349
Was, wenn ich die Daten des Clients bzw. Nutzers, der eine Datenbankanfrage bzw. -abfrage sendet, schützen muss?	350
Lösen Clean Rooms bzw. Remote Data Analysis/Access mein Datenschutzproblem?	351
Ich möchte für perfekte Privacy oder perfekte Geheimhaltung sorgen. Ist das möglich?	352
Wie stelle ich fest, ob Encrypted Computation sicher genug ist? ...	353
Wenn ich Encrypted Computation verwenden möchte, wie handhabe ich dann den Schlüsselaustausch?	354
Was ist die Privacy Sandbox von Google? Verwendet sie Encrypted Computation?	355
Data Governance und Privacy-Mechanismen	356
Warum reicht k-Anonymity nicht aus?	356
Ich denke, dass Differential Privacy nicht für meinen Anwendungsfall geeignet ist. Was kann ich stattdessen tun?	358
Kann ich mithilfe von synthetischen Daten Datenschutzprobleme lösen?	358
Wie können Daten auf verantwortungsvolle Weise weitergegeben werden, bzw. welche Alternativen gibt es zum Verkauf von Daten?	359
Wie kann ich alle privaten Informationen finden, die ich schützen muss?	360

Ich habe die persönlichen Identifikatoren entfernt, also sind die Daten jetzt geschützt, richtig?	361
Wie kann ich mit unzureichend geschützten Daten verfahren, die ich in der Vergangenheit veröffentlicht habe?	362
Ich arbeite an einem BI-Dashboard bzw. einer Visualisierung. Wie kann ich es datenschutzfreundlich gestalten?	363
Wer trifft die Entscheidungen bezüglich des Privacy Engineering? Wie kann ich Privacy Engineering in meinem Unternehmen einbinden?	364
Welche Fähigkeiten oder Vorkenntnisse benötige ich, um Privacy Engineer zu werden?	365
Warum haben Sie (Technologie oder Unternehmen hier einfügen) nicht erwähnt? Wo erhalte ich weitere Informationen? Hilfe!	366
DSGVO und Datenschutzvorschriften	367
Muss ich wirklich Differential Privacy verwenden, um Daten den Anforderungen der DSGVO/CPRA/LGPD usw. zu entziehen?	367
Ich habe gehört, dass ich personenbezogene Daten gemäß DSGVO aus berechtigtem Interesse verwenden kann. Ist das richtig?	368
Ich möchte Schrems II im Hinblick auf transatlantische Datenflüsse einhalten. Was sind mögliche Lösungen?	369
Persönliche Entscheidungen und soziale Aspekte von Privacy	370
Welche E-Mail-Provider, Browser und Anwendungen sollte ich verwenden, wenn mir meine Privatsphäre am Herzen liegt?	370
Mein Freund hat einen automatisierten Haushalts- bzw. Telefonassistenten. Ich möchte nicht, dass er mir zuhört. Was soll ich tun?	373
Ich habe mich schon lange damit abgefunden, keine Privatsphäre zu haben. Ich habe nichts zu verbergen. Warum sollte ich mich ändern?	373
Kann ich meine eigenen Daten einfach an Unternehmen verkaufen?	375
Ich mag personalisierte Werbung. Warum nicht auch Sie?	376
Hört (Füllen Sie die Lücke) gerade mit? Was kann ich dagegen tun?	377
Zusammenfassung	379
11 Machen Sie sich ans Werk und entwickeln Sie Privacy-Lösungen!	381
Überwachungskapitalismus und Data Science	381
Gig-Worker und Überwachung am Arbeitsplatz	382
Überwachung aus Gründen der »Sicherheit«	383
Luxury Surveillance	384
Massenhafte Datensammlung und Auswirkungen auf die Gesellschaft	384
Machine Learning als Datenwäsche	385
Desinformation und Fehlinformation	386

Sich zur Wehr setzen	387
Nachforschen, dokumentieren, hacken und lernen	388
Daten kollektivieren	388
Die Aufsichtsbehörden schlagen zurück	389
Die Arbeit von Communitys unterstützen	390
Als Vorkämpfer für Privacy («Privacy Champion») vorangehen	391
Ihr Privacy-Multitool	392
Vertrauenswürdige Machine-Learning-Systeme aufbauen	392
Privacy by Design	394
Privacy und Macht	396
Tschüss	398
Index	399

Angesichts der zahlreichen Vorteile der digitalen Vernetzung ist es nicht immer offensichtlich, dass futuristische Technologien auch Nachteile mit sich bringen. Instant Messaging, biometrisches Scannen, Echtzeit-Bewegungserfassung, digitaler Zahlungsverkehr und vieles mehr waren schon immer der Stoff für Science-Fiction-Fantasien. Für diejenigen unter uns, die in der Technologiebranche arbeiten (oder diese als Konsumenten erleben), ist der »Coolness-Faktor« digitaler Tools in unserer täglichen Routine schwer zu leugnen.

Die Kehrseite des digital vernetzten Lebens ist das Recht, sich vom Netz zu trennen. Für einige Tech-Millionäre der ersten Generation ist es selbstverständlich, ihre Kinder zu Hause und in der Schule vom Internet fernzuhalten. Das mag seltsam klingen, wenn man daran gewöhnt ist, die digitale Kluft als eine Trennung zwischen Besitzern von mehreren Apple-Produkten und Habenichtsen ohne 24/7-Hochgeschwindigkeitsinternet zu sehen. Da so viele unserer täglichen Interaktionen digital geworden sind, ist es jedoch für die meisten von uns eine Herausforderung, ohne unbegrenzten Onlinezugang auszukommen.

Die Nutzung digitaler Werkzeuge und der Zugang zu Onlinerräumen wird uns heute genauso angepriesen wie zu Beginn des Internets: als eine bequeme, einfache Erfahrung, die völlig freiwillig ist und Spaß macht. Aber nichts ist lustig an einer Internet-erfahrung, die sich wie ein Aufenthalt im Hotel California anfühlt – »du kannst auschecken, wann immer du willst, aber du kannst niemals abreisen«. Nichts ist fair an einer Onlinewelt, die das Offlineleben in Bezug auf alles einschränkt, was man sehen und tun kann und wie man behandelt werden könnte. Die Vorstellung, dass wir uns in der Internetwelt lediglich für eine Reihe von zwanglosen Interaktionen entscheiden, ist nicht mehr wahr: Wenn überhaupt, sind wir oft gezwungen, uns auf einer Autobahn zu bewegen, die mit Daten über uns und andere vollgestopft ist.

Viele von uns gehen fälschlicherweise davon aus, dass unsere Daten für alle anderen uninteressant sind. Aber in diesem Fall sehen wir nicht das ganze Bild. Moderne Apps und Algorithmen horten unsere Daten, um zu verknüpfen, wo wir leben, was wir verdienen, mit wem wir ausgehen und ob wir psychische Probleme oder eine sexuell übertragbare Infektion gehabt haben. Das passiert, wenn wir nicht erkennen, dass die Vorhersagefunktion von Algorithmen in der Regel dazu verwendet wird, ein

»Profil« von uns zu erstellen. Denn dafür werden Daten verwendet, die wir bereitwillig und unwissentlich zur Verfügung gestellt haben, wenn Anbieter uns Finanzprodukte, Versicherungsschutz, Arbeitsplätze, Wohnungen oder potenzielle Liebespartner verkaufen wollen (oder uns den Zugang dazu zu verwehren).

Digitale Konnektivität soll Spaß machen und sich nicht anfühlen, als würde man kriminell verfolgt. Aber genau dieses Gefühl war mein Einkaufserlebnis in der realen Welt, seit ich ein Kind in New York City war: Damals war es in der Regel alles andere als angenehm, als sichtbare Minderheit einkaufen zu gehen oder sich nach einem Taxi umzusehen. Ich kenne das Gefühl sehr gut, gescannt, überwacht und aus einer Gruppe herausgegriffen zu werden. Genau das zeigt ein Enthüllungsbericht nach dem anderen: Unsere privaten, persönlichen und dauerhaften Daten werden in »Profilen« zusammengefasst und an Datenhändler, Regierungen und Strafverfolgungsbehörden weitergegeben und zerstören somit unsere Privatsphäre. Genau wie bei verurteilten Kriminellen.

Der Schutz der Privatsphäre ist wie der Zugang zu einem Kredit oder einem guten Anwalt – etwas, das man besser hat und nicht braucht, als etwas, das man braucht und nicht hat. Es sollte nicht erst einer biometrischen Datenerfassung beim Einsteigen in ein Flugzeug bedürfen (wogegen ich kürzlich in San Francisco protestieren musste), um zu erkennen, dass unsere persönlichen Daten zu oft ohne unsere Einwilligung oder unser Wissen erhoben werden. Es sollte nicht nötig sein, dass eine Person, die einer ethnischen Minderheit angehört, einen datengesteuerten Gesundheits- oder Finanzalgorithmus als diskriminierend einstuft. Diejenigen von uns, die in der Technologiebranche tätig sind, sollten keine Gerichtsverfahren, Geldstrafen für Unternehmen oder staatliche Regulierung benötigen, um zu erkennen, dass Systeme, die unsere Daten fast zwangsweise abgreifen, uns weder Privatsphäre noch Wahlmöglichkeiten lassen. Und was ist mit denen, die ihre Privatsphäre bewahren wollen, indem sie offline bleiben? Ähnlich wie die Kreditwürdigkeit oder der Zugang zu einem guten Anwalt ist die Wahrung der Privatsphäre zum neuen Privileg der Wohlhabenden geworden.

Diese Kluft ist vielleicht das eklatanteste Problem unseres digital vernetzten Lebens. Wenn wir jemals zu einer digitalen Welt zurückkehren wollen, in die wir uns freiwillig begeben können, müssen wir den Raum begrenzen, in dem digitale Systeme ihre Fühler nach uns ausstrecken. Wenn wir den Menschen das Recht zurückgeben wollen, anonym zu surfen oder sich online zu melden, müssen wir die Mechanismen der Datenerfassung einschränken, die derzeit die meisten digitalen Systeme steuern. Mit *Data Privacy in der Praxis* bietet Frau Jarmul erprobte Techniken für den Aufbau einer Onlinewelt, die sich von der heutigen unterscheidet. Ihre Beispiele aus dem wirklichen Leben beweisen, dass man kein Privacy Engineer sein muss, um den Datenschutz sinnvoll zu gestalten.

Ich hoffe, dass alle, die sich über algorithmische Diskriminierung und »ethische Technologie« Sorgen machen, dieses Buch lesen werden. Darüber hinaus möchte ich jeden, der digitale Systeme entwirft, entwickelt oder testet, ermutigen, für sich selbst zu entscheiden, ob Datenschutz die Komponente darstellt, die unsere derzeitigen Onlineerfahrungen von denen unterscheidet, die wir wollen und brauchen.

– *Dr. Nakeema Damali Stefflbauer*
CEO, FrauenLoop and Global AI Ethics lecturer, Stanford University

Willkommen in der wunderbaren Welt des Datenschutzes! Möglicherweise haben Sie sich bereits eine Meinung zum Datenschutz (engl. *Data Privacy*) gebildet – dass er eine lästige Angelegenheit ist, dass er bürokratisch und deshalb langweilig ist, oder Sie sind vielleicht der Meinung, dass es ein Thema ist, für das lediglich Juristinnen und Juristen Interesse aufbringen können. In diesem Buch werden Sie herausfinden, wie technisch komplex und interessant die Herausforderungen des Datenschutzes sind – und auch in Zukunft sein werden. Sollte Ihre Begeisterung für knifflige mathematische und statistische Fragestellungen zu Ihrer Entscheidung geführt haben, sich mit Data Science zu befassen, dann werden Sie mit Sicherheit auch Gefallen daran finden, Datenschutz aus der Perspektive der Data Science zu erforschen. Die in diesem Buch vermittelten Inhalte werden Ihre Kenntnisse in den Bereichen Wahrscheinlichkeitstheorie, Modellierung und sogar Kryptografie erweitern.

Für Data-Science-Fachleute wird es zunehmend wichtiger, zu lernen, wie auch Datenschutzprobleme gelöst werden können. Nachdem Sie das Buch gelesen haben, werden Sie in der Lage sein, reale Probleme in Bereichen wie Cybersicherheit, Gesundheitswesen und Finanzwirtschaft zu lösen und Ihre Karriere innerhalb eines Irrgartens aus Datenschutzbestimmungen, -richtlinien und -rahmen voranzutreiben. Seit Inkrafttreten der Europäischen Datenschutz-Grundverordnung (DSGVO oder DS-GVO, engl. *General Data Protection Regulation* – GDPR) im Jahr 2018 ist die weltweite Datenschutzlandschaft noch komplexer geworden. Diese Komplexität wird weiter zunehmen, da Aufsichtsbehörden und Gesetzgeber fortwährend die Regeln dahin gehend ändern, wie, wo, warum und wann Sie Daten speichern dürfen. Wenn Sie jetzt Ihre Kompetenzen rund um den Bereich Datenschutz und Datensicherheit erweitern, ist das zweifelsohne eine sinnvolle Investition in Ihre berufliche Zukunft.

Darüber hinaus leisten Sie, wenn Sie Zeit darin investieren, neue Kenntnisse über den Datenschutz zu erlangen, einen Beitrag im Bereich der Data Science und fördern Vertrauen, Verantwortlichkeit, gegenseitiges Verständnis und soziale Verantwortung. Maschinelles Lernen (*Machine Learning*) zur Lösung von Problemen in der realen Welt stößt gegenwärtig dort auf Angst und Widerstände, wo Daten, Modelle und Systeme in nicht vertrauenswürdiger Weise genutzt wurden und sich Fragen

nach Gerechtigkeit und Fairness stellen. Ein Beispiel: Clearview AI sammelt Bilder von Gesichtern aus sozialen Netzwerken und verkauft das auf dieser Grundlage entwickelte Gesichtserkennungsmodell an Strafverfolgungsbehörden (<https://oreil.ly/PE6u1>)¹, was Fragen hinsichtlich des Eigentums an den Daten, dem Schutz der Privatsphäre und der Haftung aufwirft. Um diesem Reputationsverlust entgegenzuwirken und Wege für eine verantwortungsbewusste und vertrauenswürdige Datennutzung zu schaffen, bedarf es in der Branche Data Scientists und Machine Learning Engineers, die die vorliegenden Aufgaben und die damit verbundenen Risiken verstehen und bei der Entwicklung von Systemen diese Fragen kompetent berücksichtigen können. Der Datenschutz kann Ihnen dabei helfen, gerechtere, ethisch besser zu vertretende und verantwortungsvollere Systeme zu entwickeln, bei denen die Benutzerinnen und Benutzer die Macht und die Möglichkeit haben, sich einzubringen, und im Mittelpunkt Ihrer Ausgestaltung stehen. Mithilfe dieses Buchs können Sie diese Herausforderungen meistern und dank praxisnaher Anleitungen neue Wege finden.

Ich hoffe, dass dieses Buch einen Beitrag zur neuen Data Science leisten kann, indem es das Bewusstsein dafür schärft, wie der Schutz sensibler Daten in geeigneter Weise umgesetzt werden kann. Weltweit sind die Ängste vor der Digitalisierung persönlicher Daten – selbst für den verantwortungsvollen Einsatz durch die Regierung – so groß, dass sie die Nutzung von Daten zur Unterstützung bei sozialen Problemen wie dem Klimawandel, der Finanzaufsicht und globalen Gesundheitskrisen behindern. Wenn wir den Datenschutz in die Data Science integrieren, eröffnen sich neue Wege für die Nutzung von Daten bei wichtigen Entscheidungen für unsere Gesellschaft und unsere Welt.

Was ist Data Privacy?

Vereinfacht gesagt, schützt Privacy Daten und Menschen, indem es durch Beschränkungen hinsichtlich des Zugriffs, der Nutzung, der Verarbeitung und der Speicherung einen besseren Schutz der Privatsphäre ermöglicht und garantiert. In der Regel handelt es sich dabei um personenbezogene Daten, es umfasst aber jegliche Art der Verarbeitung. Diese Definition greift allerdings zu kurz, um Data Privacy in seiner ganzen Breite zu begreifen.

Privacy ist ein komplexes Konzept – mit Aspekten aus vielen verschiedenen Bereichen unserer Welt, sei es in rechtlicher, technischer, sozialer, kultureller oder individueller Hinsicht. Werfen wir zunächst einen Blick auf diese Aspekte und ihre Überschneidungen, damit Sie eine Vorstellung davon bekommen, wie weitreichend die Auswirkungen der in diesem Buch behandelten Themen und Vorgehensweisen sind.

¹ Eine Auflistung der URLs ohne Abkürzungen finden Sie unter <https://practicaldataprivacybook.com>.

In Abbildung E-1 sehen Sie die verschiedenen Arten der Definitionen von Privacy (bzw. des Datenschutzes oder der Wahrung der Privatsphäre)², und ich habe versucht, das jeweilige Ausmaß in der Abbildung zu illustrieren. Gehen wir sie durch und beginnen wir mit den rechtlichen Definitionen.

Im juristischen Kontext umfasst Privacy die Vorschriften, die Rechtsprechung und die Richtlinien, die festlegen, welche Maßnahmen erforderlich sind und was in einem bestimmten Staat oder einer bestimmten Gerichtsbarkeit unter Privacy zu verstehen ist. Wie Sie in den Kapiteln 1 und 8 erfahren werden, handelt es sich dabei um ein sich ständig wandelndes Rechtsverständnis und eine Landschaft, die sich in den letzten Jahren drastisch verändert hat. Es ist wichtig, dass Sie sich mit den rechtlichen Aspekten von Privacy vertraut machen, da sie sich direkt auf Ihre Arbeit auswirken können. Was passiert zum Beispiel, wenn Ihr Unternehmen von einem Audit, einer Datenschutzverletzung oder einer Verbraucherbeschwerde betroffen ist? Diese gesetzlichen Definitionen wirken sich auch auf Ihr persönliches Leben aus, beispielsweise bei der Frage, welche Rechte Sie als Datenbürger haben.

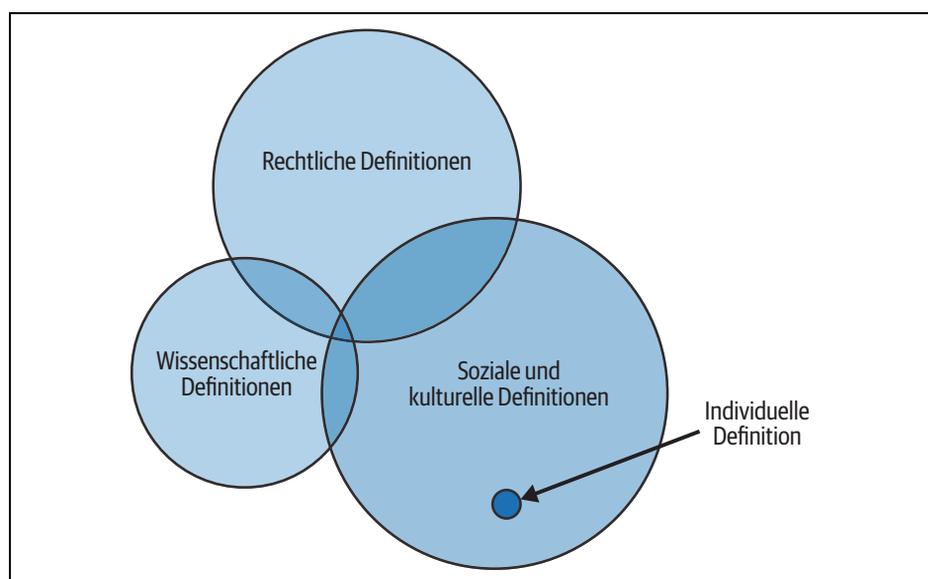


Abbildung E-1: Definitionen von Privacy

Die wissenschaftlichen bzw. technischen Definitionen von Privacy und deren Umsetzung in Ihrer täglichen Arbeit stehen im Mittelpunkt dieses Buchs. Sie lernen diese Definitionen kennen und erfahren, wie Sie wissenschaftliche Technologien zum Schutz der Privatsphäre in großem Umfang einsetzen und wie Sie technische Entscheidungen zum Thema Privacy treffen können. Mit den Tools in diesem Buch

2 Für den englischen Begriff *Privacy* gibt es im Deutschen unterschiedliche Übersetzungen. In diesem Buch wurde er je nach Kontext als *Privatsphäre* oder *Datenschutz* übersetzt. Dabei bezieht sich (Schutz oder Wahrung der) *Privatsphäre* auf den persönlichen Raum und die Freiheit vor unerwünschter Einmischung oder Überwachung. *Datenschutz* dagegen konzentriert sich eher auf den Schutz personenbezogener Daten vor Missbrauch oder unerlaubter Verwendung.

lernen Sie modernste Best Practices kennen, die in Ihrem Unternehmen möglicherweise noch nicht bekannt sind, da sie erst seit Kurzem in Produktionssystemen zur Verfügung stehen. Über diese Praktiken auf dem Laufenden zu bleiben, wird Teil Ihres Jobs sein – jedenfalls sofern Sie sich dazu entscheiden, sich auf diesen Bereich zu konzentrieren. Als technischer Experte für dieses Thema werden Sie gebeten werden, geschäftliche und juristische Entscheidungen zum Datenschutz zu unterstützen und diese in funktionsfähige Software und Systeme umzusetzen. Dies ist eine wichtige Aufgabe, denn viele der anderen Beteiligten werden kein technisches und zeitgemäßes Verständnis von Data Privacy haben.

Die sozialen und kulturellen Aspekte von Privacy lassen sich am besten anhand einer Studie zu Data Privacy von danah boyd (<http://www.danah.org>) erklären. Sie untersuchte jugendliche Mädchen und ihre Interaktion mit sozialen Medien, um zu verstehen, wie die Technologie ihr Verständnis von Konzepten wie Privacy beeinflusst. Ihre Definition lautet wie folgt:

Bei Privacy geht es weder um die Kontrolle über Daten, noch ist sie eine der Eigenschaften von Daten. Es geht um ein kollektives Verständnis der Grenzen einer gesellschaftlichen Situation und um das Wissen, wie man innerhalb dieser Grenzen agiert. Mit anderen Worten, es geht darum, die Kontrolle über eine Situation zu besitzen. Es geht darum, das jeweilige Gegenüber zu verstehen und zu wissen, wie weit Informationen verbreitet werden. Es geht darum, den Menschen, der Situation und dem Kontext zu vertrauen.

– danah boyd, in »Privacy and Publicity in the Context of Big Data«
(<https://oreil.ly/ThnPz>)

boyd weist uns mit dieser Definition auf einen neuen Aspekt von Privacy hin, der wesentliche Veränderungen bei der Gestaltung von Privacy in Systemen mit sich bringt. Im Gegensatz zu technischen und rechtlichen Definitionen stellt boyd das soziale und kulturelle Verständnis, den Kontext und die individuelle Wahl und das Bewusstsein in den Mittelpunkt. Wenn Sie ihre Arbeit lesen oder sie sprechen hören, erfahren Sie Wahrheiten, die Sie zwar oft gefühlt, aber nie vollständig erfasst haben, und zwar darüber, wie wir als Menschen und als Gesellschaft Privatsphäre und Informationen verstehen.

Wenn ich zum Beispiel meine Stimme senke und flüstere, um Ihnen etwas mitzuteilen, verstehen Sie, dass diese Information nicht für die Öffentlichkeit bestimmt ist. Wenn ich es auf einem öffentlichen Platz herausschreie und die Leute auffordere, zuzuhören, verstehen Sie, dass ich möchte, dass so viele Menschen wie möglich es hören. Wie eine Person entscheidet, mit wem sie kommuniziert, und wie sie kommuniziert, wird stark davon beeinflusst, wie diese Person den Begriff »Privacy« definiert und betrachtet (siehe Abbildung E-1). Die Fähigkeit, die eigene Kommunikation mit anderen auszuprobieren und zu verändern, hat sich im Laufe der Zeit erheblich verändert. Technologie und das Internet erlauben allen, ihre Kommunikation und die aus ihr resultierenden Möglichkeiten im Hinblick auf Privacy auf Kontexte auszuweiten, die nicht in der physischen Welt verhaftet sind. Dadurch ergeben sich neue Möglichkeiten, Kontakte zu knüpfen, sich mit anderen auszutauschen und Informationen zu teilen – und das ist wunderbar!

Diese Verlagerung von der physischen in die Onlinewelt hat jedoch auch dazu geführt, dass wir nicht mehr wissen, in welchem Kontext wir uns bewegen. Wie lauten die Regeln für diesen Raum? Wer kann mich sehen und hören? Spreche ich mit Ihnen oder mit einer Gruppe, und wie groß ist diese Gruppe? Helen Nissenbaums Forschung zur kontextuellen Integrität (<https://oreil.ly/SZ0iF>) zeigt, dass die technische Entwicklung die Wahrnehmbarkeit und Transparenz dieser Grenzen verändert hat – nicht nur über die Benutzeroberflächen, sondern auch in der grundlegenden Art und Weise, wie Systeme und Software entwickelt werden. Entscheidungen über die Standardeinstellungen von Anwendungen wirken sich auf die Privatsphäre von potenziell Millionen von Menschen gleichzeitig aus. Entscheidungen über Sicherheit und Verschlüsselung machen private Gespräche offen für Strafverfolgung und staatliche Überwachung. Data Warehouses können aus sensiblen Informationen, die nur für eine Person bestimmt sind, Zugriffsmöglichkeiten für Mitarbeitende und Datendienste Dritter schaffen. Wenn der Kontext verloren geht oder kaschiert wird und das Systemdesign die sozialen und kulturellen Definitionen von Privacy nicht berücksichtigt, hat die Technologie den menschlichen Aspekt von Privacy im Wesentlichen ignoriert.

Dieses Buch zeigt Ihnen Möglichkeiten auf, wie Sie diese gesellschaftlichen Erkenntnisse in Systemen in der Praxis umsetzen können. Sie werden viele schwierige Entscheidungen treffen müssen – aber den Nutzenden Möglichkeiten geben, sich in digitalen Räumen in Bezug auf ihre Privatsphäre zurechtzufinden; und sichere Standardeinstellungen sind Geschenke von unschätzbarem Wert, von denen die Welt mehr braucht. Während Sie dieses Buch lesen und mehr über die technischen Aspekte von Privacy erfahren, sollten Sie die soziale und die rechtliche Definition im Hinterkopf behalten – sie sind und werden für immer miteinander verwoben sein.

An wen richtet sich dieses Buch?

Dieses Buch richtet sich an Data Scientists, die sich gezielt im Bereich Data Privacy und Sicherheit weiterbilden möchten. Sie könnten dafür viele Gründe haben, wie etwa:

- Sie möchten eine Spezialisierung in Richtung »Data Privacy« verfolgen, für die Sie sich interessieren und die in der Branche eine langfristige Perspektive hat.
- Sie möchten in einen stärker regulierten Sektor wie die Finanz- oder Gesundheitsbranche wechseln, und mit diesen Kompetenzen sind Sie dort ein vielversprechender Kandidat.
- Sie arbeiten mit Forschungsdaten und würden gern eine raschere Genehmigung von Ethikkommissionen und Veröffentlichungen erhalten.
- Sie sind Freiberufler oder Berater im Bereich Data Science und möchten Ihren Kundenstamm erweitern, indem Sie kompetent mit sensiblen Daten umgehen können.
- Sie leiten ein Datenteam und möchten in der Lage sein, Produkte und Lösungen unter Berücksichtigung des Datenschutzes zu konzipieren und zu entwickeln.

- Sie möchten KI für Gutes («AI for good») einsetzen und halten den Schutz der Privatsphäre für ein wichtiges Menschenrecht.
- Ihrem Team wurde gesagt, dass Datenschutz wichtig sei, aber Sie sind sich nicht sicher, was das eigentlich bedeutet oder wie Sie es umsetzen können.
- Sie arbeiten mit sensiblen Daten und möchten sicherstellen, dass Sie sich an die Best Practices halten.
- Sie möchten ein Privacy Engineer werden und sich auf die Integration des Datenschutzes in Datenprodukte konzentrieren.
- Datenschutz und Sicherheit sind spannende Themen, und es macht Ihnen einfach Spaß, mehr darüber zu erfahren.

Ich könnte noch etliche weitere Beispiele anführen, und ich habe schon viele Menschen mit diesen unterschiedlichen Hintergründen getroffen. Eines kann ich Ihnen mit Sicherheit sagen: Die Nachfrage nach diesen Fähigkeiten steigt rapide an, und zwar nicht nur aufgrund neuer Vorschriften. Die Unternehmen investieren in diese Fähigkeiten, damit sie das Datenmanagement in eine sichere Zukunft führen können. Durch Investitionen in den Datenschutz können Unternehmen nicht nur teure Pannen vermeiden, sondern auch eine vertrauenswürdige Marke und Unternehmenskultur in Bezug auf das Datenmanagement schaffen, was sich positiv auf die Personalbeschaffung, das Marketing und die Haftung auswirkt.



Wenn Sie mit Python, Jupyter Notebooks, Mathematik und Statistik vertraut sind, werden Sie alle Abschnitte gut verstehen können. Sie können diesen tiefergehenden theoretischen und implementierungsorientierten Abschnitten folgen, aber bei der Lektüre auch weglassen, solange Sie die grundlegenden Konzepte verstehen.

Machen Sie sich keine Sorgen, wenn Sie sich schon länger nicht mehr mit Mathematik beschäftigt haben. Zu jedem der Beispiele habe ich Ihnen eine Erklärung mitgeliefert. Es wird Ihnen helfen, sich beim Durchlesen Zeit zu lassen.

Beim Schreiben dieses Buchs habe ich Feedback von Softwareentwicklerinnen und -entwicklern, Sicherheitsspezialisten und sogar Datenschutzanwälten erhalten, denen dieses Buch nützlich erschien. Obwohl diese Leute nicht meine Zielgruppe sind, hoffe ich, dass dieses Buch jedem helfen kann, der sich für Privacy und Technologie sowie deren Überschneidung in Datenystemen interessiert.

Beim Lesen dieses Buchs und beim Durcharbeiten der Übungen werden Sie sehen, wie Aspekte der Data Privacy die Wunder der Data Science hervorheben, die Sie bereits kennen und lieben. Wie in anderen herausfordernden Bereichen der Data Science führt Sie dieses Buch von einfachen Methoden für die Lösung im Bereich Privacy zu schwierigeren Methoden, von denen einige noch nicht vollständig gelöst sind. Genau wie bei der linearen Regression, die »einfach funktioniert«, möchten Sie mit einfachen und offensichtlichen Lösungen beginnen. Aber wenn die Lösung, die Sie benötigen, über die einfache Lösung hinausgeht, müssen Sie detailliertere Fragen stellen, die technische und ethische Implikationen haben. Diese Fragen zu finden

und sie und ihre Antworten zu erforschen, wird Sie zu einem besseren Data Scientist und Technologen oder einer besseren Statistikerin und Mathematikerin machen.

Vielleicht ist dieses Buch alles, was Sie benötigen, um ein Technologie zu werden, der über ein paar zusätzliche Kenntnisse im Bereich Data Privacy verfügt. Das ist okay! Vielleicht ist dieses Buch aber auch das erste von mehreren Büchern, das Sie weiter in dieses Gebiet führt. Sollte das für Sie verlockend klingen, möchte ich Sie nun mit dem Konzept des Privacy Engineering vertraut machen.

Privacy Engineering

Ich gehe davon aus, dass der Bereich Privacy Engineering (<https://oreil.ly/XENvQ>) in den nächsten zehn Jahren weiter an Bedeutung gewinnen wird.³ Die Fähigkeiten, die Sie in diesem Buch erwerben, indem Sie die Übungen durcharbeiten und das neu erlangte Wissen auf Ihre Arbeit anwenden, werden Sie auf diese Rolle vorbereiten.

In Unternehmen, in denen Data Science ein wichtiger Bestandteil ist, ist ein Privacy Engineer zum Teil Data Scientist und zum Teil Engineer. Das bedeutet, dass Sie im Gegensatz zu anderen Rollen in der Data Science aktiv an der Entwicklung und Architektur von Lösungen arbeiten, anstatt Daten zu untersuchen oder eine Idee in einer experimentellen Umgebung zu testen. Das könnte bedeuten, dass Sie direkt mit den Data-Engineering-Teams, den Software- bzw. Anwendungsteams oder sogar den Systemarchitektinnen Ihres Unternehmens zusammenarbeiten, um sicherzustellen, dass Data Privacy sowohl in den Produkten als auch in den internen Anwendungen berücksichtigt wird. Dies gilt für alle Datenströme von Verbrauchern und Mitarbeitenden, für Software, die für das Datenmanagement verwendet wird, sowie für interne und externe Datenverwendungszwecke. Im Rahmen dieser Arbeit müssen Sie die Grundlagen der Technik und der Architektur verstehen, insbesondere was die Entwicklung von Systemen und die Integration von Systemen untereinander betrifft. Zu diesen Themen gibt es einige verwandte Bücher, mit denen Sie sich befassen können:

- *Software Architecture in Practice*, 4th Edition (<https://oreil.ly/5M2Zt>)
- *Handbuch moderner Softwarearchitektur* (<https://dpunkt.de/produkt/handbuch-moderner-softwarearchitektur/>)
- *Entwurfsmuster von Kopf bis Fuß* (<https://dpunkt.de/produkt/entwurfsmuster-von-kopf-bis-fuss-2/>)
- *Datenintensive Anwendungen designen* (<https://dpunkt.de/produkt/datenintensive-anwendungen-designen/>)
- *Practical MLOps* (<https://oreil.ly/tXioO>)

Um bestmöglich zu bestimmen, welche Tools und welche Software für ein Unternehmen geeignet sind, ist eine ausgeklügelte Architektur erforderlich. Die einfache Implementierung von Datenschutzrichtlinien durch Plug-and-play-Anbieter greift

³ Anmerkung: In der Regel vermeide ich es, Vorhersagen zu treffen, da man oft falsch liegt. Diese hier basiert jedoch auf den Erfahrungen, die ich in den letzten sechs Jahren in der Branche gewonnen habe.

daher oft zu kurz, um diese Probleme zu lösen. Abgesehen davon bedeutet die wachsende Zahl von Anbietern von Datenschutztechnologien, dass Sie zum Entscheidungsträger werden, wenn es darum geht, Technologien zu entwickeln oder zu kaufen und für das Datenschutzmanagement einzusetzen. Dabei werden Sie die in diesem Buch gelernten Konzepte anwenden, um Bewertungskriterien aufzustellen, Fragen zur Implementierung zu stellen und die Flexibilität, den Support und die Produktmerkmale zu analysieren. In dieser Rolle werden Sie feststellen, wie gut potenzielle Anbieter die Anforderungen Ihres Unternehmens erfüllen können, da die Abhängigkeit von privaten, sensiblen und vertraulichen Daten wächst.

Ein Privacy Engineer ist nicht einfach nur ein weiterer Data Scientist oder Data Architect, der sich um die Einhaltung des Datenschutzes sorgt, letztlich aber keine Befugnis, keine Zeit und kein Budget zur Verfügung hat, Entscheidungen bezüglich Data Privacy treffen zu können. Es ist zwar erfreulich, dass das Engagement (engl. *Advocacy*) Teil der Rolle des Data Scientist geworden ist, aber beim Privacy Engineering geht es darum, Privacy-Techniken zu entwickeln und diese anzuwenden, wenn Daten eingespeist (engl. *ingest*), gesammelt, transformiert, gespeichert und schließlich in Data-Science-Anwendungen eingesetzt werden. Das Eintreten für Privacy mag vielleicht hilfreich sein, aber erst die Umsetzung erbringt den Beweis, dass diese Technologien funktionieren.

Ein Privacy Engineer ist auch nicht nur ein Data Engineer, der sich mit Privacy beschäftigt. Privacy Engineers können zwar mit Data Engineers zusammenarbeiten – und werden oft für ein Projekt oder ein Proof of Concept in ein Team eingegliedert –, aber sie müssen mit verschiedenen Teilen des Unternehmens zusammenarbeiten und werden in viele Projekte einbezogen, bei denen ihr Fachwissen gefragt ist. Als Spezialistinnen und Spezialisten sind sie nicht allzu lange an ein einzelnes Projekt oder einen Anwendungsfall gebunden. Ihr Wissen ist vielmehr eine ungeheuer wertvolle Ressource, die für die dringendsten geschäftlichen Fragestellungen im Zusammenhang mit Data Privacy eingesetzt werden sollte.

Das Berufsbild des Privacy Engineer ist noch nicht ausdefiniert begriffen und erfährt eine stetige Weiterentwicklung. Obwohl größere Technologieunternehmen mittlerweile aktiv Personal für diese Position einstellen, erinnert mich das Aufkommen dieser Berufsbezeichnung an das Aufkommen des Begriffs *Machine Learning Engineer* im Jahr 2018. Privacy Engineering – also der Umgang mit dem Datenschutz in der Praxis – ist eine relativ neue Qualifikation im Bereich Data Science, die sich aufgrund der Bedürfnisse und Anforderungen der Branche entwickelt. Ich bin gespannt, wie sich die Rolle des Privacy Engineer in zwei oder auch in zehn Jahren darstellen wird –, und hoffe, dass dieses Buch dazu beiträgt, ein paar weitere Menschen für diesen Bereich zu begeistern.

Warum ich dieses Buch geschrieben habe

Als das Thema Data Privacy für mich zum ersten Mal interessant wurde, kam es mir wie ein riesiges Labyrinth vor. Der Großteil der Materialien war für mich nicht ver-

ständig, und die einführenden Leitfäden wurden oft von Menschen geschrieben, die mir ihre Software verkaufen wollten. Glücklicherweise kannte ich ein paar Leute in der Data-Privacy-Community, die mir dabei halfen, ein tieferes und umfassenderes Verständnis erlangen zu können. Es bedurfte vieler Stunden des Studiums und zahlreicher hilfsbereiter Personen, damit ich mich von einem neugierigen Data Scientist zu jemandem entwickeln konnte, der die Themen, die Sie in diesem Buch antreffen, beherrscht. Ich kann Ihnen verraten, ich lerne weiterhin jedes Jahr aufs Neue dazu und tauche tiefer in das Gebiet ein.

Ich bin davon überzeugt, dass die Fähigkeiten, die Sie in diesem Buch erlernen werden, heute und auch künftig für Data Scientists unerlässlich sind. Meine eigene Lernkurve verlief viel zu steil. Und genau das soll Ihnen dieses Buch ersparen. Ich habe dieses Buch geschrieben, um Ihnen eine ansprechende, schnellelebige und praxisorientierte Umgebung zu verschaffen, in der Sie dazulernen, Fragen stellen, hilfreiche Ratschläge finden und sich näher mit den anspruchsvollen Themen befassen können.

Dieses Buch ist als ein nützlicher Überblick gedacht, der Ihnen dabei hilft, den Datenschutz ohne Vorkenntnisse aktiv in Ihre Arbeit zu integrieren. Sie lernen gängige Strategien wie Pseudonymisierungs- und Anonymisierungsverfahren und neuere Ansätze wie Berechnungen auf Basis verschlüsselter Daten (*Encrypted Computation*) und Federated Data Science kennen. Wenn dieses Buch als Sprungbrett für Ihre akademische Karriere dient oder dazu verhilft, dass Sie als Forscherin tätig werden, wäre das großartig. Das Berufsfeld braucht intelligente und neugierige Menschen, die an ungelösten Problemen in diesem Bereich arbeiten möchten. Doch im Großen und Ganzen ist dieses Buch ein praxisorientierter Überblick, der, sollten Sie mehr wissen wollen, unterwegs Verweise liefert.

Data Scientists und Technologen, die Datenschutz- und Sicherheitsthemen in ihre tägliche Arbeit miteinbeziehen müssen, werden dieses Buch hilfreich finden. Es gibt einige Kapitel, die Ihnen als Kurzreferenz dienen, während Sie durch die Welt der Data Privacy navigieren. Wenn Sie das Buch von Anfang bis Ende lesen, werden Sie eine solide Kenntnis über die Materie erlangen und lernen, wie Sie neue, Ihnen zuvor unbekannte Datenschutzprobleme lösen können. Ein kurzes Nachschlagen liefert Ihnen unkomplizierte Ratschläge dazu, wie Sie mit bestimmten Datenschutznotfällen umgehen können, die in Ihrer täglichen Arbeit auftauchen.

Aufbau des Buchs

Dieses Buch soll Ihnen einen praktischen Ansatz für Data Privacy bieten und enthält eine Mischung aus Theorie, Übungen und Anwendungsfällen. Dabei gliedert es sich in die folgenden Kapitel:

- In Kapitel 1, *Data Governance und einfache Datenschutzansätze*, geht es um Data Governance und einfache Ansätze zur Data Privacy. Sie erhalten eine Reihe von Hinweisen zur Verwaltung von Daten, zur Rückverfolgung von Einwilligungen und zur Pseudonymisierung von Daten, die Sie intern verwenden möchten.

- Mit Kapitel 2, *Anonymisierung*, tauchen Sie in das Thema Anonymisierung ein und erfahren, welche modernen Ansätze Sie heute verwenden können und wie das US Census Bureau Differential Privacy als Werkzeug für Data Scientists entwickelt hat.
- In Kapitel 3, *Datenschutz in Datenpipelines integrieren*, erfahren Sie, wie Sie damit beginnen können, Data Privacy in Datenpipelines und -workflows zu automatisieren, wobei verschiedene Anwendungsfälle rund um die Themen Einwilligung (engl. *Consent*), Anonymisierung (engl. *Anonymization*) und Data Engineering aufgezeigt werden.
- Kapitel 4, *Angriffe auf die Privatsphäre*, gibt Ihnen einen Überblick über die bisher bekannten Angriffe auf die Privatsphäre, z.B. wie der Netflix-Price-Datensatz de-anonymisiert wurde, und zeigt Ihnen, wie Sie mögliche Sicherheitslücken und Angriffe bei der Arbeit mit sensiblen Daten erkennen können.
- Kapitel 5, *Machine Learning und Data Science datenschutzkonform gestalten*, befasst sich damit, wie Machine Learning datenschutzkonform gestaltet werden kann und wie Sie Datenschutz in Data-Science-Projekte integrieren können. Dieses Kapitel sollte als Schnellreferenz verwendet werden, um bestimmte Ansätze in einem projekt- oder produktbezogenen Data-Science-Team zu evaluieren.
- Kapitel 6, *Federated Learning und Data Science*, erläutert, wie föderale Ansätze beim Machine Learning (Federated Learning) und in der Data Science funktionieren und vergleicht diese mit anderen Privacy-Ansätzen und Datensparsamkeit.
- In Kapitel 7, *Encrypted Computation*, finden Sie Informationen zum Thema Encrypted Learning und Encrypted Computation für Data Privacy in der Data Science, wobei Sie sich mit Multiparty Computing und homomorphen Verschlüsselungsprotokollen (engl. *Homomorphic Encryption Protocols*) und -bibliotheken beschäftigen.
- Kapitel 8, *Datenschutzrechtliche Aspekte*, vermittelt Ihnen, wie Sie Datenschutzbestimmungen und -richtlinien interpretieren und anwenden können. Dabei werden die DSGVO, das kalifornische Verbraucherschutzgesetz (California Consumer Privacy Act, CCPA) und verschiedene Beispiele für interne Richtlinien vorgestellt, die Sie dabei unterstützen, die rechtliche Seite von Privacy zu durchdringen.
- Kapitel 9, *Datenschutz und Anwendungen aus der Praxis*, hilft Ihnen, das Gelernte anzuwenden, um sichere und private Datensysteme in realen Anwendungsfällen zu konzipieren. Dieses Kapitel dient ebenfalls als Schnellreferenz, insbesondere für Data Architects und das Data-Science-Management.
- Kapitel 10, *Häufig gestellte Fragen und ihre Antworten!*, fasst häufig gestellte Fragen und Anwendungsfälle zusammen und dient dementsprechend als praktische Referenz für Datenschutznotfälle. So können Sie selbstbewusst vorangehen und sicherstellen, dass Data Privacy in jedem Projekt und in Ihrem norma-

len Arbeitsablauf integraler Bestandteil ist. Außerdem erfahren Sie mehr über die sozialen und auch persönlichen Aspekte von Privacy und können diese auf Ihr Privatleben übertragen.

- Das letzte Kapitel des Buchs, Kapitel 11, *Machen Sie sich ans Werk und entwickeln Sie Privacy-Lösungen!*, soll Ihnen helfen und Sie dazu motivieren, Ihre neu erworbenen Datenschutzkenntnisse dafür zu nutzen, das Fachgebiet und Ihren eigenen Weg weiter voranzutreiben!

Die im englischsprachigen Buch enthaltenen Links wurden der Einfachheit halber zu O'Reilly-URLs verkürzt. Diese URLs unterliegen nur einem Mindestmaß an Tracking und wurden auf Konformität mit der DSGVO und auf den Schutz der Privatsphäre überprüft. Sollte Ihnen dieses Maß an Tracking nicht zusagen, können Sie die vollständige Liste der URLs des englischsprachigen Buchs unter <https://practicaldataprivacybook.com> einsehen.

Neuerungen in der deutschsprachigen Ausgabe

Die deutsche Übersetzung dieses Buchs enthält einige zusätzliche Abschnitte und überarbeitete Passagen, um der zunehmenden Verbreitung von Large Language Models (LLMs) und GPT-basierten Anwendungen Rechnung zu tragen. Diese Ergänzungen sollen das Bewusstsein für Angriffe auf die Privatsphäre bei der Verwendung dieser Modelle schärfen und den aktuellen Stand der Technik in Bezug auf den Schutz bzw. die Bereitstellung datenschutzfreundlicher generativer KI-Dienste aufzeigen.

In diesem Buch verwendete Konventionen

Die folgenden typografischen Konventionen werden in diesem Buch verwendet:

Kursiv

Kennzeichnet neue Begriffe, URLs, E-Mail-Adressen, Dateinamen und Dateinamen.

Konstante Zeichenbreite

Wird für Programmlistings und für Programmelemente in Textabschnitten wie Namen von Variablen und Funktionen, Datenbanken, Datentypen, Umgebungsvariablen, Anweisungen und Schlüsselwörter verwendet.

Konstante Zeichenbreite, fett

Kennzeichnet Befehle oder anderen Text, den der Nutzer wörtlich eingeben sollte.

Konstante Zeichenbreite, kursiv

Kennzeichnet Text, den der Nutzer je nach Kontext durch entsprechende Werte ersetzen soll.



Dieses Symbol steht für einen Tipp oder eine Empfehlung.



Dieses Symbol steht für einen allgemeinen Hinweis.



Dieses Symbol warnt oder mahnt zur Vorsicht.

Verwenden von Codebeispielen

Zusätzliche Materialien (Codebeispiele, Übungen und so weiter) können Sie unter <https://github.com/kjam/practical-data-privacy> herunterladen.

Wir haben eine Webseite für dieses Buch, auf der wir Errata, Beispiele und zusätzliche Informationen veröffentlichen. Sie können diese Seite unter <https://oreil.ly/practicalDataPrivacy> aufrufen.

Dieses Buch dient dazu, Ihnen bei der Erledigung Ihrer Arbeit zu helfen. Im Allgemeinen dürfen Sie die Codebeispiele aus diesem Buch in Ihren eigenen Programmen und der dazugehörigen Dokumentation verwenden. Sie müssen uns dazu nicht um Erlaubnis bitten, solange Sie nicht einen beträchtlichen Teil des Codes reproduzieren. Beispielsweise benötigen Sie keine Erlaubnis, um ein Programm zu schreiben, in dem mehrere Codefragmente aus diesem Buch vorkommen. Wollen Sie dagegen eine CD-ROM mit Beispielen aus Büchern von O'Reilly verkaufen oder verbreiten, benötigen Sie eine Erlaubnis. Eine Frage zu beantworten, indem Sie aus diesem Buch zitieren und ein Codebeispiel wiedergeben, benötigt keine Erlaubnis. Eine beträchtliche Menge Beispielcode aus diesem Buch in die Dokumentation Ihres Produkts aufzunehmen, bedarf hingegen unserer ausdrücklichen Zustimmung.

Wir freuen uns über Zitate, verlangen diese aber nicht. Ein Zitat enthält Titel, Autor, Verlag und ISBN. Beispiel: »*Data Privacy in der Praxis* von Katharine Jarmul, O'Reilly 2024, ISBN 978-3-96009-233-9.«

Wenn Sie glauben, dass Ihre Verwendung von Codebeispielen über die übliche Nutzung hinausgeht oder außerhalb der oben vorgestellten Nutzungsbedingungen liegt, kontaktieren Sie uns bitte unter komentar@oreilly.de.