

STUDIE DES VERBUNDPROJEKTS »CYBERSICHERHEIT FÜR DIE DIGITALE VERWALTUNG«

## PRIVACY UND BIG DATA



---

## Impressum

---

### Redaktion

Anna Spiegel

### Layout und Satz

Matthias Buss

### Kontakt

Fraunhofer-Institut für  
Sichere Informationstechnologie SIT  
Rheinstraße 75  
64295, Darmstadt

### Bildquellen

Seite XI: Hessisches Ministerium des Innern und für Sport  
Seite 97: Bild erstellt von Martin Steinebach und Christian Winter, lizenziert unter CC BY-SA 4.0; abgeleitet von »Road in Jökulsá á Dal« von floheinstein auf Flickr, lizenziert unter CC BY-SA 2.0.

© Fraunhofer-Institut für  
Sichere Informationstechnologie SIT,  
Darmstadt, 2020

---

## Hinweise

---

Die in diesem Dokument enthaltenen Arbeitsergebnisse sind sorgfältig und unter Zugrundelegung des bekannten Standes der Wissenschaft erstellt worden, stellen jedoch Forschungsansätze dar. Eine Haftung oder Garantie dafür, dass die Arbeitsergebnisse bzw. Informationen die Vorgaben der aktuellen Rechtslage erfüllen, wird aus diesem Grund nicht übernommen. Gleiches gilt für die Brauchbarkeit, Vollständigkeit oder Fehlerfreiheit, so dass jede Haftung für Schäden ausgeschlossen wird, die aus der Benutzung dieser Arbeitsergebnisse bzw. Informationen entstehen können. Diese Haftungsbeschränkung gilt nicht in Fällen von Vorsatz.

Alle Bezeichnungen für Personen, die in dieser Studie genannt werden, gelten sowohl für das männliche als auch das weibliche Geschlecht. Der Einfachheit halber wird durchgehend das generische Maskulinum verwendet.

# PRIVACY UND BIG DATA

## STUDIE

---

Verbundprojekt »Cybersicherheit für die digitale Verwaltung«

---

Christian Winter  
Martin Steinebach  
Wendy Heereman  
Simone Steiner  
Verena Battis  
Oren Halvani  
York Yannikos  
Christoph Schübler

Teilweise basierend auf Arbeiten von  
Marcel Schäfer, Anika Pflug, Jamal Pasha

16. September 2020

# INHALTSVERZEICHNIS

Inhaltsverzeichnis . . . . .	IV
Vorwort . . . . .	VIII
Zusammenfassung . . . . .	IX

---

## **I Einführung** **11**

---

1 Einleitung . . . . .	12
2 Risiken und Herausforderungen für den Datenschutz durch Big Data. . . . .	13
3 Scheitern von Privatheit . . . . .	15
3.1 NYC Taxi . . . . .	15
3.2 AOL . . . . .	16
3.3 Netflix . . . . .	16
3.4 FindFace . . . . .	17
Literatur . . . . .	19

---

## **II Datenschutzrechtliche Perspektive** **20**

---

Abkürzungsverzeichnis . . . . .	21
4 Begriffsbestimmungen . . . . .	22
5 Problemstellung . . . . .	24
6 Abgrenzung personenbezogener, pseudonymisierter und anonymisierter Daten nach DSGVO. . . . .	25
6.1 Anonymisierung . . . . .	26
6.2 Pseudonymisierung . . . . .	27
6.3 Datenschutzrechtliche Relevanz von Big-Data-Anwendungen unter Verarbeitung anonymisierter Daten . . . . .	27
7 Regelungen der DSGVO und des neuen Bundesdatenschutzgesetzes im Überblick	28
7.1 Erlaubnistatbestände der DSGVO . . . . .	28
7.2 Grundsätze des Datenschutzes . . . . .	30
7.2.1 Verbotsprinzip . . . . .	30
7.2.2 Zweckbindung . . . . .	30
7.2.3 Zweckänderung . . . . .	30
7.2.4 Zeitliche Speicherbegrenzung . . . . .	31

7.2.5	Datenminimierung	32
7.2.6	Transparenz	32
7.2.7	Datenrichtigkeit	32
7.2.8	Integrität und Vertraulichkeit	32
7.2.9	Rechenschaftspflicht	33
7.2.10	Einwilligung zur Datenverarbeitung im Rahmen von Big Data	33
7.3	Rechte der betroffenen Person	33
7.3.1	Recht auf Information	34
7.3.2	Auskunftsrecht	35
7.3.3	Berichtigung	36
7.3.4	Löschung	36
7.3.5	Einschränkung der Verarbeitung	37
7.3.6	Recht auf Datenübertragbarkeit	37
7.3.7	Widerspruch	37
7.3.8	Sonstige Rechte der betroffenen Person	38
7.3.9	Profiling	38
7.4	Technischer und organisatorischer Datenschutz	39
7.4.1	Gemeinsame Bedingungen der Art. 24, 25, 32 DSGVO	39
7.4.2	Sicherstellung der Rechtskonformität durch technische und organisatorische Maßnahmen – Artikel 24 DSGVO	41
7.4.3	Art. 25 Datenschutz durch Technikgestaltung und durch datenschutzfreundliche Voreinstellungen	41
7.4.4	Art. 32 – Gewährleistung der Datensicherheit	43
7.4.5	Standard-Datenschutzmodell (SDM) und Leitlinien der Artikel-29-Datenschutzgruppe zur Orientierung	44
7.5	Datenschutz-Folgenabschätzung	45
7.6	Verzeichnis von Verarbeitungstätigkeiten	47
7.7	Auftragsverarbeitung	48
7.8	Übermittlung an Drittländer oder an internationale Organisationen	49
<hr/>		
<b>8</b>	<b>Empfehlungen für Big-Data-Anwender</b>	<b>50</b>
8.1	Privacy by Design und by Default	50
8.2	Regelmäßige Prüfung der Anonymisierungsverfahren sowie sonstigen technischen und organisatorischen Maßnahmen	51
8.3	Mitarbeiterschulung	51
8.4	Interne Datenschutzrichtlinie	52
8.5	Durchführung einer Datenschutz-Folgenabschätzung	52
8.6	Anforderungen an IT-Systeme	53
8.7	Genehmigte Verhaltensregeln und Zertifizierung der Big-Data-Anwendungen gemäß DSGVO	53
<hr/>		
<b>9</b>	<b>Zusammenfassung</b>	<b>54</b>
	Literatur	55

---

## III Technische Ansätze zum Schutz der Privatsphäre 57

---

<b>10</b>	<b>Schutzziele bei Big Data</b>	<b>58</b>
10.1	Gründe für die Absicherung von Daten und Anforderungen in Big-Data-Szenarien	58
10.1.1	Schutz der Privatsphäre	58
10.1.2	Wissenschutz	59
10.2	Anforderungen für einen sicheren und datenschutzgerechten Umgang mit Daten	59
10.2.1	Übertragung (»Data in Transit«)	59
10.2.2	Speicherung (»Data at Rest«)	60
10.2.3	Verarbeitung (»Data in Use«)	60
10.3	Ausblick auf die Beiträge dieses Teils der Studie	61
<hr/>		
<b>11</b>	<b>Privacy by Design und Big Data</b>	<b>62</b>
11.1	Konzepte	63
11.2	Von der Leitlinie zur technischen Umsetzung	64
11.3	Diskussion	66
<hr/>		
<b>12</b>	<b>Verschlüsselungsmechanismen für Daten in Big-Data-Systemen</b>	<b>67</b>
12.1	Grundbegriffe	67
12.2	Betriebsmodi für Blockchiffren	69
12.3	Zustände digitaler Daten in Big-Data-Systemen	70
12.3.1	Data at Rest	70
12.3.2	Data in Transit	70
12.3.3	Data in Use	70
12.4	Geeignete Verschlüsselung für »Data at Rest«	71
12.5	Geeignete Verschlüsselung für »Data in Transit«	71
<hr/>		
<b>13</b>	<b>Geeignete Verschlüsselung für Data in Use</b>	<b>74</b>
13.1	Homomorphe Verschlüsselung	74
13.1.1	Privatsphären-Homomorphismen	75
13.1.2	Partiell homomorphe Verschlüsselung (PHE)	76
13.1.3	»Begrenzt« homomorphe Verschlüsselung (SHE)	77
13.2	Sichere Mehrparteienberechnung	78
<hr/>		
<b>14</b>	<b>Anonymisierung strukturierter Daten</b>	<b>80</b>
14.1	Anonymisierung auf Basis von k-Anonymität	81
14.1.1	k-Anonymität	81
14.1.2	l-Diversität	82
14.1.3	t-Nähe	83
14.1.4	Algorithmen für k-Anonymität und verwandte Kriterien	83
14.2	Differential Privacy	84
14.2.1	Differential Privacy für Frage-Antwort-Systeme	85
14.2.2	Weitere Anwendungen für Differential Privacy	86

<b>15 Anonymisierung von Texten</b>	<b>87</b>
15.1 Anonymisierung auf der Metadatenebene	87
15.2 Anonymisierung auf der Inhaltsebene	87
15.3 Anonymisierung auf der Schreibstilebene	89
<b>16 Anonymisierung im Kontext von maschinellem Lernen</b>	<b>90</b>
16.1 Privatsphärenrisiken beim maschinellen Lernen	90
16.1.1 Risiken durch gelernte Modelle	90
16.1.2 Risiken durch maschinelles Lernen in der Cloud	91
16.2 Privatsphärenfreundliches maschinelles Lernen	91
16.2.1 Differential Privacy für maschinelles Lernen	92
16.2.2 Maschinelles Lernen mit homomorpher Verschlüsselung	92
16.2.3 Kollaboratives maschinelles Lernen	92
<b>17 Zu bewältigende Herausforderungen</b>	<b>94</b>
17.1 Technische Fortschritte	94
17.1.1 Effektive und effiziente Anonymisierung	95
17.1.2 Berechnungen auf verschlüsselten Daten	96
17.1.3 Synthetische Daten	97
17.2 Zusammenspiel von gesetzlichen Anforderungen und Technologie	98
17.2.1 Konkretisierung erforderlicher Maßnahmen und Handlungen	98
17.2.2 Definition von Personenbezug und Anonymität	99
17.2.3 Verhindern von Diskriminierung	101
17.3 Unterstützung der involvierten Parteien	101
17.3.1 Anwendbarkeit von Privacy by Design	101
17.3.2 Schulung von Praktikern	101
17.3.3 Finanzielle Förderung von Forschung und Entwicklung	102
17.3.4 Erleichterung des Zugangs zum Datenschutz für Endnutzer	102
17.4 Langfristige Herausforderungen	102
17.4.1 Effiziente voll-homomorphe Verschlüsselung	102
17.4.2 Beweisbare Anonymisierung	103
17.4.3 Rechtliche Harmonisierung	104
17.4.4 Nicht absehbare Entwicklungen der Gesellschaft	104
17.5 Fazit	105
<b>18 Publikationen</b>	<b>CVI</b>
Literatur	CVII

# VORWORT

Digitalisierung ist und bleibt eine der drängendsten Herausforderungen unserer Gesellschaft. Besonders die Corona-Pandemie hat uns seit der ersten Hälfte des Jahres 2020 gezeigt, wie unabhkmmlich und notwendig Digitalisierung und innovative Technologien sind. Mit der rasanten Entwicklung der Digitalisierung nimmt die Bedeutung von Sicherheitsaspekten stetig zu – gerade bei großen Datenmengen. Die Analyse solcher Datenmengen mithilfe von künstlicher Intelligenz nimmt in vielen Anwendungen eine Schlüsselrolle ein, denn solche Big-Data-Systeme können oft als Hebel für Innovationen dienen. Big Data kann den Alltag vieler Menschen erleichtern und dabei helfen, die Herausforderungen in vielen Bereichen wie beispielsweise Verkehr, Infrastruktur und Medizin besser zu bewältigen.

Gleichzeitig haben wir eine vielseitige Debatte über die Vereinbarkeit von Datenschutz und Datensicherheit mit der Nutzbarmachung von personenbezogenen Daten für einen gesamtgesellschaftlichen Nutzen geführt. Der Schutz des Privatlebens und personenbezogener Daten sind europäische Grundrechte, die wiederum Grundvoraussetzung für demokratische Werte und die Ausübung anderer Grundrechte sind. Das weltweite Datenaufkommen wird jeden Tag um ein Vielfaches größer. Umso wichtiger ist es, erforderliche Maßnahmen zum Schutz der Privatsphäre zu treffen – insbesondere in sensiblen Bereichen wie Gesundheitsdaten. „Privacy by Design“ – also das Einbinden von Schutzaspekten hinsichtlich der Privatheit von Beginn an – muss fester Bestandteil bei Big-Data-Projekten sein. Die Sicherheit der Bürgerinnen und Bürger muss auch in diesem Kontext stets ein hohes Gut sein.

Die Hessische Landesregierung verfolgt beim Thema Schutz in der digitalen Welt mit dem Verbundprojekt „Cybersicherheit für die digitale Verwaltung“ einen ganzheitlich orientierten Ansatz, u.a. mithilfe der Förderung der Cybersicherheitsforschung durch das Hessische Ministerium des Innern und für Sport. Die vorliegende Studie „Privacy und Big Data“ ist ein Ergebnis dieser Forschungsförderung. In dieser Studie werden die Herausforderungen von Big Data in Verbindung mit dem Datenschutz und dem Umgang mit personenbezogenen Daten analysiert und Handlungsempfehlungen abgeleitet, die von großem Nutzen sind. Mit dieser Studie wollen wir einen Beitrag dazu leisten, die Wichtigkeit des Datenschutzes im Kontext von Big Data und Privatsphärenschutz aufzuzeigen und den Weg für innovative digitale Lösungen, auch innerhalb der Verwaltung auf Landes- und Kommunalebene, zu bereiten.

## **Peter Beuth**

Hessischer Minister des Innern und für Sport



*Peter Beuth*



# ZUSAMMENFASSUNG

Die vorliegende Studie adressiert das Thema Privatheit im Kontext von Big Data. Sie besteht aus drei Hauptteilen: In der **Einführung** werden abstrakt die Risiken und Herausforderungen für den Datenschutz durch Big Data beschrieben. Zum Abschluss der Einführung zeigen wir einige Beispiele, die Datenschutzrisiken im Kontext von Big Data demonstrieren. Dabei werden nicht Angriffe auf und Sicherheitslücken von Systemen betrachtet, sondern Risiken für die Privatsphäre, die durch heutige Analysemöglichkeiten entstehen. Diese Analysemöglichkeiten erlauben vermeintlich anonymisierte Daten zu deanonymisieren oder öffentliche Daten aus ganz unterschiedlichen Kontexten auf privatsphärenfeindliche Weise miteinander zu verknüpfen. Der zweite Teil adressiert die **Datenschutzrechtliche Perspektive**. Hier werden die rechtlichen Rahmenbedingungen des Datenschutzes systematisch erläutert und allgemeine Ratschläge für technische und organisatorische Maßnahmen zum Datenschutz gegeben. Im dritten Teil adressieren wir **Technische Ansätze zum Schutz der Privatsphäre**. Hier werden verschiedene technischen Möglichkeiten, insbesondere kryptographische Methoden und Anonymisierungsansätze, beschrieben.

Durch Big Data ist die Bedeutung von **Daten als Ressource** deutlich geworden. Besonders der Wert und der Schutz **personenbezogener Daten** wird dabei immer wieder diskutiert, da heute das Sammeln detaillierter Profile und dadurch detaillierte Einblicke in das Verhalten von Personen möglich sind.

Sollen personenbezogene Daten verfügbar gemacht werden, sollte davor eine **Anonymisierung** erfolgen, der Personenbezug also nicht mehr herstellbar sein. Allerdings liefert Big Data auch Werkzeuge, um genau diese Vorkehrungen wieder aufzuheben. Durch das Verknüpfen unterschiedlicher Datenquellen konnte in der Vergangenheit Anonymisierung erfolgreich gebrochen werden, beispielsweise im Fall des Netflix Prize Datensatzes.

Aus rechtlicher Perspektive ist Privatheit in Big Data natürlich in erster Linie ein Thema, das durch die **DSGVO** adressiert wird. Hier ist zuerst notwendig, Daten, die sich auf identifizierte oder identifizierbare natürliche Person beziehen und somit personenbezogene Daten sind, von nicht personenbezogenen Daten zu unterscheiden. Personenbezogene Daten sind wenn möglich zu **anonymisieren** oder zu **pseudonymisieren**. Dabei ist zu beachten, dass künftige Entwicklungen die Schutzmechanismen schwächen können und daher ihr Wirksamkeit regelmäßig geprüft werden sollte.

Sollen personenbezogene Daten verarbeitet werden, die nicht anonymisiert werden, so gelten zahlreiche **datenschutzrechtliche Anforderungen**. Zuerst muss dabei geprüft werden, ob eine Verarbeitung auf Basis eines Erlaubnistatbestandes zulässig ist. Ist dies der Fall, gelten die Grundsätze des Datenschutzes, also Zweckbindung, Zeitliche Speicherbegrenzung, Datenminimierung, Datenrichtigkeit, Integrität und Vertraulichkeit, Rechenschaftspflicht und Regeln zur Einwilligung. Außerdem hat die betroffene Person zahlreiche Rechte, wie das Informationsrecht, Auskunftsrecht, das Recht zur Berichtigung, eine Einschränkung der Verarbeitung, das Recht auf Datenübertragbarkeit und ein Widerspruchsrecht.

Um gesetzeskonform in BigDataUmgebungen personenbezogene Daten verarbeiten zu können, empfehlen sich eine Reihe von Vorgehensweisen. Orientierung bieten hier beispielsweise **Privacy by Design und by Default**, die zahlreiche organisatorische und technische Impulse beinhalten. Die Wirksamkeit der daraus abgeleiteten Maßnahmen sollte regelmäßig geprüft werden. Aufgrund der Komplexität und Bedeutung der Fragestellung empfehlen sich gezielte **Mitarbeiterschulungen** und das Aufstellen von **Datenschutzrichtlinien**. Um Bedeutung und Risiko präzise abschätzen zu können, ist eine **Datenschutz-Folgenabschätzung** ein etabliertes Werkzeug.

**Verschlüsselung** ist auch in BigDataUmgebungen eine wichtige Basistechnologie. Sie wird insbesondere auf gespeicherte Daten («at rest») und bei den Datenübertragung («in transit») eingesetzt. Ansätze, Daten im verschlüsselten Zustand zu verarbeiten, sind ebenfalls bekannt, aber deutlich komplexer und weniger verbreitet. Hierunter fallen die **homomorphe Verschlüsselung** und die **sichere Mehrparteienberechnung**.

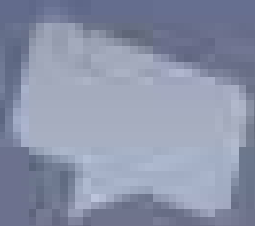
Die Anonymisierung von Daten kann grob in die Betrachtung strukturierter Daten (wie Datenbanken) und unstrukturierter Daten (wie Texten) unterschieden werden. Bei strukturierten Daten sind dabei verschiedenen Strategien wie Generalisierung, Löschung, Mikroaggregation und Verfälschung bekannt. Bekannte Ansätze sind dabei **kAnonymität** und **Differential Privacy**.

Sollen **Texte anonymisiert** werden, ist die erste Herausforderung zu erkennen, wo personenbezogene Daten vorhanden sind. Dies kann auf der Metadatenebene, der Inhaltsebene und der Schreibstilebene der Fall sein. Während Metadaten mit wenig Aufwand anonymisiert werden können, erfordert dies auf der Inhaltsebene ein zuverlässiges Erkennen von allen Worten mit Personenbezug wie Namen, aber auch anderen Hinweisen, die eine Person identifizieren können. Diese Stellen können dann auf verschiedene Weise unkenntlich gemacht oder ersetzt werden. Soll auch der Schreibstil keinen Hinweis auf einen Autor geben, dann ist **Author Obfuscation** notwendig, wodurch eine stilistische Normalisierung erfolgen soll.

Big Data und **maschinelles Lernen** werden heute häufig gemeinsam eingesetzt. Dadurch entstehen neue Fragestellungen der Privatheit. So konnten Risiken durch gelernte Modelle nachgewiesen werden, da diese auf bekannte Daten stärker reagieren als auf unbekannte, ähnliche Daten. Allerdings wird auf die Risiken reagiert, das privatsphärenfreundliche maschinelle Lernen ist hier ein Ansatz. Es beinhaltet Konzepte wie Differential Privacy für maschinelles Lernen, maschinelles Lernen mit homomorpher Verschlüsselung und kollaboratives maschinelles Lernen.

Am Ende der Studie finden sich **zahlreiche Herausforderungen**, die das Themenfeld Privacy und Big Data noch bewältigen muss, um als gelöst betrachtet werden zu können. Diese umfassen notwendige technische Verbesserungen, aber auch organisatorische Ansätze wie Schulungen und einen vereinfachten Zugang von Betroffenen zu Mechanismen des Datenschutzes oder notwendige rechtliche Harmonisierungen.

# I EINFÜHRUNG



# 1 Einleitung

In den letzten Jahren sind Daten zu einer sehr wertvollen Ressource geworden. Big Data ist der Oberbegriff für die jüngste Entwicklung in den Bereichen Datenerfassung, Datenverarbeitung und Analysetools. Kombiniert mit den bloßen Möglichkeiten der Datenerfassung über viele mit dem Internet verbundenen Geräten, führt diese Entwicklung zu erheblichen Veränderungen im sozialen (digitalen) Leben. Sie stellt neue Möglichkeiten, aber auch neue, weitreichende Herausforderungen dar.

Datenschutz und Big-Data-Analysen auf personenbezogenen Daten sind ein Thema, über das regelmäßig diskutiert wird und welches zumeist widersprüchlich zu sein scheint. Datenschutz erfordert eine Vielzahl von Einschränkungen und im Falle der Umsetzung von Privacy-by-Design (PbD) einen deutlichen Mehraufwand beim Systementwurf. Big-Data-Analysen, insbesondere auch bei der Abwehr von terroristischen Bedrohungen auf der anderen Seite erfordern schnellen und pragmatischen Zugriff auf möglichst viele Datenquellen, um schnell ein Lagebild zu erstellen oder eine Bedrohungslage abzuschätzen. Überspitzt gesehen führt dies zu einer Entscheidung zwischen rechtskonformem oder ergebnisorientiertem Handeln.

Heute sind allerdings bereits eine Reihe technischer Ansätze bekannt, die Datenschutz und Datenanalyse leichter vereinbar erscheinen lassen und so die oben genannten Widersprüche reduzieren helfen: Daten können unter Betrachtung ihrer weiteren Nutzbarkeit anonymisiert werden, um so einen Kompromiss zwischen Anonymität und Aussagekraft zu erreichen.

Systeme können so gestaltet werden, dass Datenzugriffe zwar ermöglicht, aber auch zuverlässig protokolliert werden, um einen Missbrauch der Daten nachvollziehbar zu gestalten. Viele weitere Bausteine sind bekannt, um ein Big-Data-System datenschutzfreundlich zu gestalten.

## 2 Risiken und Herausforderungen für den Datenschutz durch Big Data

In diesem Abschnitt identifizieren wir eine Reihe von Risiken, die sich aus der Verwendung von Big-Data-Systemen zur Verarbeitung personenbezogener Daten ergeben. Die Spezialität von Big-Data-Systemen ist die Verarbeitung großer Datenmengen (»Volumen«, engl. »Volume«) aus verschiedenen Quellen unterschiedlicher Art und in heterogenen Formaten und Strukturen (»Vielfalt«, engl. Variety«) und die Bereitstellung von Ergebnissen in kurzer Zeit (»Geschwindigkeit«, engl. »Velocity«). Daher sind Volume, Variety und Velocity die drei klassischen »V«s, die Big Data charakterisieren. Einige Autoren haben auch ein oder mehrere zusätzliche »V«s hinzugefügt. Wir glauben jedoch, dass die Kombination der drei klassischen »V«s bereits erfasst, was Big-Data-Systeme von traditionellen Datenverarbeitungssystemen unterscheidet, während die anderen »V«s nicht so gut geeignet sind, zwischen Big-Data-Systemen und anderen Systemen zu unterscheiden.

Basierend auf den Eigenschaften von Big-Data-Systemen existieren die folgenden Bedrohungen für Einzelpersonen:

### 1. Sammlung detaillierter Profile:

Traditionelle Systeme wurden für einen bestimmten Zweck gebaut, und diese Systeme hatten eine klar definierte Datenquelle oder zumindest eine klar definierte Art von Datenquellen. Darüber hinaus war die Art der gesammelten Daten begrenzt und gut definiert. Im Gegensatz dazu zielt Big Data darauf ab, viele Eingaben aus mehreren Quellen und idealerweise aus mehreren Domänen zusammenzufassen. So birgt Big Data das Risiko, detaillierte und facettenreiche Profile von Personen zu erstellen. Man betrachte zum Beispiel die Tracking-Möglichkeiten auf Internetseiten und in Smartphone-Apps. Die Beobachtung von Endverbrauchern mit solchen Methoden ergibt oft einen breiten Querschnitt ihres privaten und beruflichen Lebens.

### 2. Mehr Einblicke in Individuen:

Big Data verwendet ausgefeilte analytische Techniken, um ansonsten verborgene Zusammenhänge zu entdecken. Die Philosophie von Big Data besteht nicht darin, eine bestimmte Frage zu beantworten, sondern eine Vielzahl von Mustern zu entdecken, die für die Beantwortung verschiedener bestehender Fragen nützlich sein können

oder die für die Beantwortung zukünftiger Fragen relevant sein könnten. Darüber hinaus versucht Big Data oft, das zukünftige Verhalten von Personen vorherzusagen. Daher besteht die Gefahr eines breiten und tiefen Einblicks in das Privatleben der Menschen.

### 3. Verknüpfung vermeintlich nicht personenbezogener Daten mit Personen:

Ein wichtiger Ansatz, um Datenschutzrisiken zu reduzieren und Beschränkungen durch Datenschutzgesetze zu umgehen, ist die Anonymisierung von Daten. Korrekt anonymisierte Daten sind nicht mehr personenbezogen, da sie keine spezifischen Informationen über einzelne Personen mehr enthalten. Wahre Anonymisierung ist jedoch schwierig, und die Geschichte ist voll von Beispielen für gescheiterte Anonymisierungen, z. B. beim Datenschutz des Netflix Prize (s. Abschnitt 3.3). So bleibt immer das Risiko einer erneuten Identifizierung von Personen in vermeintlich anonymisierten Daten bestehen. Dies ist zwar kein reines Problem von Big Data, aber die Analysekapazitäten von Big-Data-Systemen bergen das Risiko, dieses Problem zu verstärken. Solche Systeme könnten versteckte Zusammenhänge in angeblich anonymisierten Daten verwenden, um Erkenntnisse über Personen zu gewinnen, ohne diese Zusammenhänge explizit zu machen. Eine weitere Variante dieses Problems ergibt sich, wenn Informationen aus einer ursprünglich nicht personenbezogenen Domäne mit Personen verknüpft werden können. Dies geschieht, wenn die für die Verknüpfung relevanten Attribute ausreichend feingranular werden. Eine solche Art von Attributen sind etwa geographische Koordinaten. Beispiele für Informationen, die mit Personen anhand von geographischen Koordinaten verknüpft werden können, sind visuelles Material, das von modernen Kartendiensten bereitgestellt wird, sowie feingranulare Wetterdaten wie Blitzaufzeichnungen.

Der Ansatz von Big Data ist konträr zu den Prinzipien des Datenschutzes. Der Datenschutz beschränkt die Menge der zu erhebenden Daten, den Umfang der Verarbeitung und die Dauer der Speicherung. Big Data strebt jedoch das Gegenteil an: Sammle so viele Daten wie möglich, versuche, jede erdenkliche Einsicht zu erhalten, und lösche die Daten niemals.

Daher steht Big Data im Konflikt mit dem Datenschutzrecht, wenn es um personenbezogene Daten geht, was sehr häufig vorkommt. Das bedeutet, dass Big-Data-Systeme und -ansätze beim Umgang mit personenbezogenen Daten angepasst werden müssen. Betreiber solcher Systeme und – noch mehr – diejenigen, die für die Verarbeitung verantwortlich sind, müssen die datenschutzrechtlichen Anforderungen wie Zweckbindung, Datenminimierung und Speicherzeitbegrenzung einhalten (vgl. Abschnitt 7.2).

Neben den Maßnahmen, die von Anwendern von Big Data verlangt werden, um die Datenschutzgesetze einzuhalten, stellt Big Data auch die verfügbaren Technologien zum Datenschutz sowie die bestehenden Datenschutzgesetze wie die DSGVO, vor neue Herausforderungen. Insbesondere die folgenden Schwierigkeiten ergeben sich aus der Natur von Big Data:

Anonymisierungstechniken müssen gegen Risiken der Offenlegung gestärkt werden, die im Zeitalter von Big Data deutlich zugenommen haben. Darüber hinaus muss die Anonymisierung effizient sein, um anwendbar zu sein. Zusätzlich müssen Anonymisierungstechniken so viele Informationen wie möglich von den Originaldaten erhalten. Andernfalls ist die Anonymisierung für die meisten Big-Data-Anwendungen nutzlos.

Die Möglichkeit, personenbezogene Informationen aus vermeintlich anonymisierten oder ursprünglich nicht personenbezogenen Eingabedaten zu extrahieren, verwischt die Grenze zwischen personenbezogenen und nicht personenbezogenen Daten. Daher müssen bei der Definition personenbezogener und nicht personenbezogener Daten und bei den Regeln für den Umgang mit solchen Daten diese in der Praxis auftretenden Probleme berücksichtigt werden.

Die Analyse personenbezogener Daten in großen Daten-systemen gibt Anlass zur Sorge über Diskriminierung. Dieses Problem besteht bereits bei traditionellen Scoringverfahren (z. B. Kreditwürdigkeitsprüfung oder Scoring in Bezug auf (zukünftige) kriminelle Affinität). Dieses Problem verschärft sich jedoch im Zusammenhang mit großen Datenmengen. Solche Systeme verwenden komplexe Entscheidungsalgorithmen – oft basierend auf maschinellem Lernen –, die selbst für die Entwickler der Systeme undurchsichtig sind. Daher können solche Systeme verzerrte Trainingsdaten oder unvorsichtige ausgewählte Attribute in systematische Diskriminierung verwandeln.

Machine-Learning-Algorithmen, die in Big-Data-Systemen verwendet werden, erfordern geeignete und aufbereitete Trainings- und Testdaten. Solche Daten sind jedoch aus mehreren Gründen schwer zu erhalten. Im Zusammenhang mit personenbezogenen Daten ist der Datenschutz ein wesentlicher Grund für die Schwierigkeit, Schulungs- und Testdaten zu erhalten. Dies gilt für Forschung und Entwicklung im akademischen und im wirtschaftlichen Umfeld. Anonymisierte Daten sind eine Möglichkeit, dieses Problem zu lösen. Synthetische Daten sind eine weitere Option.

Eine große Herausforderung für datenschutzfreundliche Lösungen ist der Wettbewerb mit Lösungen, die sich nicht so sehr um den Datenschutz kümmern. Häufig haben die Betreiber der letztgenannten Lösungen einen Marktvorteil durch die Nutzung der erfassten Daten. Darüber hinaus entscheiden sich die Endnutzer eher für die Lösung mit der größten Nützlichkeit, gemessen an der Einfachheit der Bedienung und der gewünschten Funktionalität – ohne Rücksicht auf den Datenschutz (und einen Großteil der erweiterten Funktionalität). Daher ist es eine entscheidende Herausforderung, die Endnutzer zu erreichen, wenn man datenschutzfreundliche Lösungen anbietet.

Der Stand der Technik zur Behandlung dieser Herausforderungen wird in Teil III beleuchtet. Die Limitationen des Stands der Technik und die Aufgaben zur Lösung dieser Limitation werden in Kapitel 17 herausgearbeitet.

## 3 Scheitern von Privatheit

In diesem Kapitel fassen wir Vorfälle zusammen, die aufzeigen, wie Privatheit und Big Data in Konflikt geraten können. Dabei geht es oft darum, wie verfügbare Daten mehrerer

Quellen zusammengefügt werden können, um vermeintlich anonyme oder anonymisierte Daten realen Personen zuzuordnen.

### 3.1 NYC Taxi

Im März 2014 veröffentlichte die New York City Taxi and Limousine Commission<sup>1</sup> Taxi-Reiseprotokolle von mehr als 173 Millionen Reisen, um eine Visualisierung und Analyse zu erstellen. Abbildung 3.1 zeigt beispielhaft einen Datensatz aus diesem Bestand. Der primäre Zweck dieser Daten war es, die Forschungsgemeinschaft zu unterstützen und herauszufinden, welche Teile der Stadt unterversorgt waren und zu erfahren, wie hoch die Wahrscheinlichkeit ist, um 4 Uhr morgens ein Taxi zu bekommen. Diese Protokolle enthielten den Ort, wo die Taxifahrt begann und wo sie endete, Datum und Uhrzeit der Taxifahrt, wie viel Trinkgeld für die Taxifahrt bezahlt wurde und anonymisierte Hinweise auf die Identitäten von Taxi und Fahrer. Allerdings hat die New York City Taxi and Limousine Commission die MD5-Hash-Funktion verwendet, um die Lizenz- und Medaillonnummern zu hashen, um die Offenlegung der Identität und die Privatsphäre zu verhindern. Um die tatsächliche Struktur des Datensatzes für die Visualisierung und Analyse beizubehalten, haben jedes Taxi und jeder Fahrer denselben MD5-Hash im anonymisierten Datensatz. Pandurangan [6] zeigte, dass diese Logs deanonymisiert werden können. Er stellte dar, dass in New York eine Taxi-Lizenz eine 6-stellige Nummer oder eine 7-stellige Nummer ist, die mit der Nummer 5 beginnt, was 3 Millionen unterschiedlichen Taxi-Lizenz-Nummern ergibt. Die Medaillon-Nummern hat ein Muster aus Ziffern und Buchstaben, was zu 19 Millionen möglichen Medaillon-Nummern führt. Er zeigte, dass man durch die Berechnung von 22 Millionen MD5 Hashes und dem Abgleich dieser berechneten MD5 Hashes mit dem veröffentlichten Datensatz den gesamten Datensatz von Taxi-Logs komplett de-anonymisieren kann. Durch die Verwendung des de-anonymisierten Datensatzes kann man unter anderem die Einnahmen von Taxifahrern ermitteln. Darüber hinaus können mit den Lizenznummern und Medaillonnummern weitere persönliche Informationen von Taxifahrern offengelegt wer-

den. Weiterhin kann man mit einigen zusätzlichen öffentlich zugänglichen Informationen herausfinden, wer in der Taxifahrt gereist ist, indem man den veröffentlichten Datensatz auf Google Maps visualisiert.

**Datensatz:** 6B111958A39B24140C973B262EA9FEA5,  
D3B035A03C8A34DA17488129DA581EE7, VTS, 5, ,  
2013-12-03 15:46:00, 2013-12-03 16:47:00, 1, 3660,  
22.71, -73.813927, 40.698135, -74.093307, 40.829346  
**Feldbezeichnungen:** medallion, hack\_license, vendor\_id,  
rate\_code, store\_and\_fwd\_flag, pickup\_datetime,  
dropoff\_datetime, passenger\_count,  
trip\_time\_in\_secs

Abbildung 3.1: Beispiel für einen Datensatz und Bezeichnungen der Felder.

Quelle: Pandurangan [6].

So fand man zum Beispiel heraus, dass berühmte Persönlichkeiten wie Bradley Cooper und Olivia Munn auf Reisen waren, als sie das Taxi nahmen, von wo aus sie das Taxi nahmen und sogar, wie viel sie bezahlten.

<sup>1</sup> <https://www1.nyc.gov/site/tlc/index.page>.

## 3.2 AOL

America Online (AOL) war eine der bekanntesten Suchmaschinen Amerikas. Am 3. August 2006 veröffentlichte AOL einen Datensatz, der Suchanfrage-Datensätze von 658.000 Personen über den Zeitraum von 3 Monaten enthält [5]. Der Zweck dieser Datenfreigabe war es, Forschungsarbeiten zu initiieren, die darauf abzielen, menschliches Verhalten und Interessen aus den Suchanfragen heraus zu untersuchen. Es war ein Goldschatz für Data-Mining-Analysten und -Forscher aus verschiedenen Communities, da nur selten zuverlässige und qualitativ hochwertige Daten wie Suchanfragen veröffentlicht werden. Von der Mehrheit der Suchmaschinen werden Informationen über Suchanfragen der Nutzer streng vertraulich behandelt.

Um Datensätze für die Forschungsgemeinschaft nutzbar zu machen und die Identität und Privatsphäre der betroffenen Personen zu schützen, hat AOL die expliziten Identifikatoren wie IP-Adressen und Benutzernamen durch künstliche, eindeutige Identifikatoren, d. h. neu generierte Pseudonyme, ersetzt. Diese Pseudonyme halfen den Forschern, die verschiedenen Suchanfragen den entsprechenden Benutzern zuzu-

ordnen. Es war notwendig, dass jede betroffene Person im gesamten Datensatz dieselbe Kennung hat, damit Verhaltensanalytiker das Verhalten und die Interessen der betroffenen Person über den Zeitraum hinweg analysieren können.

Die AOL-Datenveröffentlichung löste eine heftige Diskussion unter den Forschern über den Schutz der Privatsphäre der betroffenen Personen aus. Einige Forscher versuchten herauszufinden, wie viel Privatsphäre und Informationen offengelegt werden können – entweder durch eine erneute Identifizierung der betroffenen Person oder durch die Lokalisierung der betroffenen Person mithilfe von Suchanfragen.

Es dauerte nicht lange, bis zwei Journalisten der New York Times eine Witwe mithilfe ihrer Suchanfragen ausfindig machten und identifizierten [5]. Dies war eine starke Verletzung der Privatsphäre, weil der Rest ihrer Recherchen ein bemerkenswert detailliertes und trauriges Bild ihres Lebens gemalt hat. AOL wurde heftig für eine derart massive Offenlegung von Informationen und Datenschutzverletzungen kritisiert. Daraufhin entfernte AOL schnell den gesamten Datensatz.

## 3.3 Netflix

Einige Monate nach dem Scheitern der AOL-Datenfreigabe, startete Netflix, einer der bekanntesten Online-Unterhaltungsdienstleister am 2. Oktober 2006 den Netflix Prize<sup>2</sup>. Dies war ein Wettbewerb mit einem Preisgeld von einer Million Dollar, um den Algorithmus für Filmempfehlungen, basierend auf vorherigen Filmbewertungen und gesehenen Filmen, zu verbessern.

Alle Teilnehmer erhielten Zugang zu einer Trainingsmenge von 100 Millionen Datensätzen und einer Testmenge von 1,4 Millionen Bewertungen von 480.000 zufällig ausgewählten anonymisierten Nutzern zu 18.000 Filmen. Jeder Datensatz enthielt Filmtitel, Filmbewertung und Bewertungsdatum. Netflix entfernte explizite Identifikatoren der Nutzer aus dem Datensatz und ersetzte sie durch eindeutige Identifikatoren, um den Bezug zu den individuellen Nutzern für die Wettbewerbsteilnehmer aufrechtzuerhalten.

Im Gegensatz zu AOL hat Netflix den Datensatz nicht primär für die allgemeine Forschung freigegeben, sondern für die Unterstützung eigener wirtschaftlicher Zwecke. Netflix wollte seine Popularität und seinen Umsatz steigern, indem es die Qualität seiner Empfehlungen für Filme und Serien verbesserte. Wie AOL wurde auch der Netflix-Datensatz von Forschern und Informatikern sehr geschätzt als eine Möglichkeit, ihre Ideen und Algorithmen auf echten und hochwertigen Daten zu entwickeln und auszuwerten.

Netflix behauptete, die Privatsphäre und die persönlichen Daten der Betroffenen zu schützen. Innerhalb von nur zwei Wochen nach Beginn des Wettkampfes zeigten jedoch Narayanan und Shmatikov [4], dass in dem veröffentlichten Datensatz von Netflix fast alle Nutzer mit sehr wenig Zusatzinformation wiedererkannt werden können.

Mithilfe von Wahrscheinlichkeitsrechnung demonstrierten

<sup>2</sup> <https://www.netflixprize.com/rules.html>.



Narayanan und Shmatikov, wie viele Informationen ein Angreifer über Netflix-Abonnenten benötigt, um sie im Trainingsdatensatz zu identifizieren und schließlich die gesamten Informationen über die Abonnenten zu erhalten. Sie gingen davon aus, dass ein Angreifer von acht Filmen weiß, die ein Nutzer bewertet hat, dass er das Bewertungsdatum auf zwei Wochen genau kennt und dass er bei mindestens sechs dieser Filme die korrekte Bewertungszahl kennt. Anhand dieser Annahmen zeigten sie, dass 99 % der Netflix-Abonnenten im Trainingsdatensatz individuell identifizierbar sind. Nimmt man an, dass der Angreifer als Hintergrundwissen nur zwei Bewertungen mit dem Datum auf drei Tage genau und mit korrekter Bewertungszahl kennt, so können immerhin 68 % der Nutzer identifiziert werden und für die restlichen Nutzer

ist die Anzahl der Kandidaten sehr gering.

Narayanan und Shmatikov demonstrierten mithilfe von Nutzerbewertungen auf IMDb<sup>3</sup>, die frei verfügbar und offen ist, dass das nötige Hintergrundwissen oft leicht zu erlangen ist und solche Angriffe in der Praxis erfolgreich sind. Mit der Annahme, dass die Nutzer, die sowohl auf IMDb und Netflix aktiv sind, oft denselben Film, den sie auf IMDb bewerten, auch auf Netflix bewerten, nahmen eine Stichprobe von 50 IMDb-Nutzern und verglichen sie mit dem Trainingsdatensatz von Netflix, um zu zeigen, dass IMDb-Daten ausreichen, um Abonnenten in Netflix-Datensätzen zu identifizieren. In der Tat konnten sie so zwei Netflix-Nutzer mit extrem hoher Sicherheit identifizieren.

### 3.4 FindFace

FindFace<sup>4</sup> ist eine App, die es erlaubt, Fotos von Personen zu erstellen und diese dann in sozialen Medien zu suchen. Ermöglicht wird dies mit Algorithmen zur Gesichtserkennung der Firma ntech lab<sup>5</sup> und dem sozialen Netzwerk VK<sup>6</sup>.

Die Erkennungsrate bei der Lösung ist hoch; in einem Test von Kaspersky [1] wurden bei gut geeigneten Aufnahmen Raten von 90 % erreicht. Natürlich müssen die Personen, die erkannt werden soll, auch im sozialen Netzwerk mit einem Foto hinterlegt sein. Die Zahl von Personen, auf die dies zutrifft, wird allerdings immer größer. Weiterhin ist die Auswahl von VK nur willkürlich. Die gleiche Vorgehensweise ist auch auf Facebook, LinkedIn oder XING möglich, wenn die entsprechenden Fotos und Daten der Nutzer gecrawlt werden.

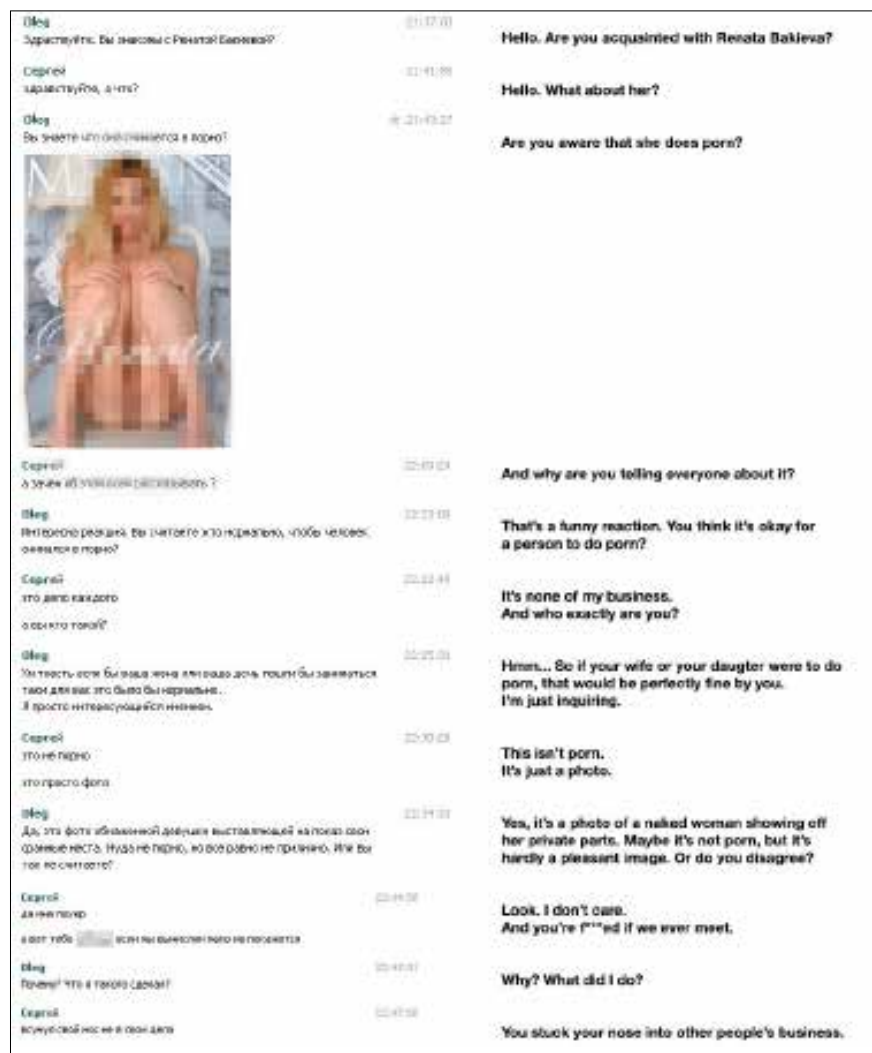


Abbildung 3.2: Screenshot eines Dialogs zwischen einem Dvach-Nutzer und einem Bekannten einer Pornodarstellerin. Übersetzung durch GlobalVoices. Quelle: GlobalVoices [7].

3 <https://www.imdb.com/>.

4 [findface.ru](http://findface.ru).

5 [ntechlab.com](http://ntechlab.com).

6 [vk.com](http://vk.com).

FindFace wurde vom NIST [3] evaluiert. Dort wurde bei einer vorgegebenen Falsch-Positiv-Rate von 0,001 eine Falsch-Negativ-Rate von 0,22 erreicht. Jede 5. Person wird also korrekt erkannt, wenn sie sich in der Referenzdatenbank befindet, wenn akzeptiert wird, dass bei jeder 1000. Überprüfung ein Fehlalarm ausgelöst wird.

### **Dvach**

Nutzer des Forums Dvach verwendeten FindFace, um Porno-Darstellerinnen und Prostituierte zu identifizieren und dann die sozialen Kontakte der Frauen mit deren Aktivitäten zu konfrontieren [7]. Abbildung 3.2 zeigt ein Beispiel eines solchen Dialogs. Angeblich geschah diese Aktion der De-Anonymisierung und Ansprache als Kritik an der Porno-Industrie. Der Vorfall gilt als drastisches Beispiel für die Risiken, die die Kombination aus Gesichtserkennung und sozialen Medien mit sich bringt.

### **Kommerzialisierung**

Die Erotik-Industrie nutzt inzwischen Gesichtserkennung zur Unterstützung von Suchvorgängen durch Anwender.

### **Nutzung gegen Demonstranten**

Es existieren inzwischen zahlreiche Hinweise, dass FindFace inzwischen in Russland dazu eingesetzt wird, Teilnehmer von Demonstrationen zu identifizieren. Dazu werden hochauflösende Bilder während der Demonstration erstellt und die Gesichter der Teilnehmer dann mit der mit der App abgeglichen. Dies geschieht teilweise durch anonyme Dritte, die die Teilnehmer denunzieren, es gibt aber angeblich auch Festnahmen durch die Polizei [2].

## Literatur

- [1] Vladislav Biryukov. *Ihr Gesicht kann nicht ausgetauscht werden*. 22. Apr. 2016. URL: <https://www.kaspersky.de/blog/findface-experiment/7505/>.
- [2] Conrad Conrad. *FindFace: Verbrecherjagd mit neuester Gesichtserkennungssoftware*. 26. Juli 2017. URL: <https://www.datenschutz-notizen.de/findface-verbrecherjagd-mit-neuester-gesichtserkennungssoftware-0118560/>.
- [3] Patrick Grother u. a. *The 2017 IARPA Face Recognition Prize Challenge (FRPC)*. NIST Interagency Report (NISTIR) 8197. National Institute of Standards und Technology (NIST), Nov. 2017. DOI: 10.6028/NIST.IR.8197. URL: <https://www.nist.gov/programs-projects/face-recognition-prize-challenge-2017>.
- [4] Arvind Narayanan und Vitaly Shmatikov. »Robust De-anonymization of Large Sparse Datasets«. In: *Proceedings of the 2008 IEEE Symposium on Security and Privacy*. IEEE Computer Society, Mai 2008, S. 111–125. ISBN: 978-0-7695-3168-7. DOI: 10.1109/SP.2008.33.
- [5] Paul Ohm. »Broken Promises of Privacy: Responding to the Surprising Failure of Anonymization«. In: *UCLA Law Review* 57.6 (Aug. 2010). Hrsg. von Darcy Pottle, S. 1701–1777. URL: <http://uclalawreview.org/pdf/57-6-3.pdf>.
- [6] Vijay Pandurangan. *On Taxis and Rainbows – Lessons from NYC’s improperly anonymized taxi logs*. 21. Juni 2014. URL: <https://tech.vijayp.ca/of-taxis-and-rainbows-f6bc289679a1>.
- [7] Kevin Rothrock. *Facial Recognition Service Becomes a Weapon Against Russian Porn Actresses*. 22. Apr. 2016. URL: <https://advox.globalvoices.org/2016/04/22/facial-recognition-service-becomes-a-weapon-against-russian-porn-actresses/>.

## II DATENSCHUTZRECHTLICHE PERSPEKTIVE



## Abkürzungsverzeichnis

Abs.	Absatz
Art.	Artikel
BDSG	Gesetz zur Anpassung des Datenschutzrechts an die Verordnung (EU) 2016/679 und zur Umsetzung der Richtlinie (EU) 2016/680 (Bundesdatenschutzgesetz)
DSGVO	Datenschutz-Grundverordnung
DSFA	Datenschutz-Folgenabschätzung
Erwgr.	Erwägungsgrund
EWR	Europäischer Wirtschaftsraum
lit.	Buchstabe (lateinisch litera)
Rn.	Randnotiz
TOM	Technische und organisatorische Maßnahmen

## 4 Begriffsbestimmungen

Bei den nachfolgenden Begriffsbestimmungen handelt es sich um eine Zusammenstellung der für diese Arbeit wichtigsten Legaldefinitionen des Art. 4 DSGVO.

### Personenbezogene Daten

Alle Informationen, die sich auf eine identifizierbare natürliche Person (betroffene Person) beziehen; als identifizierbar wird eine natürliche Person angesehen, die direkt oder indirekt, insbesondere mittels Zuordnung zu einer Kennung wie einem Namen, zu einer Kennnummer, zu Standortdaten, zu einer Online-Kennung oder zu einem oder mehreren besonderen Merkmalen identifiziert werden kann, die Ausdruck der physischen, physiologischen, genetischen, psychischen, wirtschaftlichen, kulturellen oder sozialen Identität dieser natürlichen Person sind.

### Besondere Kategorien personenbezogener Daten

Daten aus denen die rassische und ethnische Herkunft, politische Meinungen, religiöse oder weltanschauliche Überzeugungen oder die Gewerkschaftszugehörigkeit hervorgehen, sowie die Verarbeitung von genetischen Daten, biometrischen Daten zur eindeutigen Identifizierung einer natürlichen Person, Gesundheitsdaten oder Daten zum Sexualleben oder der sexuellen Orientierung einer natürlichen Person.

### Verarbeitung

Jeder mit oder ohne Hilfe automatisierter Verfahren ausgeführte Vorgang oder jede solche Vorgangsreihe im Zusammenhang mit personenbezogenen Daten wie zum Beispiel das Erheben, das Erfassen, das Ordnen, die Speicherung, das Auslesen, das Abfragen, die Verwendung, das Löschen oder die Vernichtung.

### Einschränkung der Verarbeitung

Ist die Markierung gespeicherter personenbezogener Daten mit dem Ziel, ihre künftige Verarbeitung einzuschränken.

### Pseudonymisierte Daten

Personenbezogene Daten, die in einer Weise verarbeitet werden, dass die personenbezogenen Daten ohne Hinzuziehung zusätzlicher Informationen nicht mehr einer spezifischen

betroffenen Person zugeordnet werden können, sofern diese zusätzlichen Informationen gesondert aufbewahrt werden und technischen und organisatorischen Maßnahmen unterliegen, die gewährleisten, dass die personenbezogenen Daten nicht einer identifizierten oder identifizierbaren natürlichen Person zugewiesen werden.

### Anonymisierte Daten<sup>1</sup>

Daten, die sich nicht auf eine identifizierte oder identifizierbare natürliche Person beziehen, oder personenbezogene Daten, die in einer Weise anonymisiert worden sind, dass die betroffene Person nicht oder nicht mehr identifiziert werden kann.

### Einwilligung

Jede freiwillig für den bestimmten Fall, in informierter Weise und unmissverständlich abgegebene Willensbekundung in Form einer Erklärung oder einer sonstigen eindeutigen bestätigenden Handlung, mit der die betroffene Person zu verstehen gibt, dass sie mit der Verarbeitung der sie betreffenden personenbezogenen Daten einverstanden ist.

### Verantwortlicher

Die natürliche oder juristische Person, Behörde, Einrichtung oder andere Stelle, die allein oder gemeinsam mit anderen über die Zwecke und Mittel der Verarbeitung von personenbezogenen Daten entscheidet.

### Auftragsverarbeiter

Eine natürliche oder juristische Person, Behörde, Einrichtung oder andere Stelle, die personenbezogene Daten im Auftrag des Verantwortlichen verarbeitet.

### Empfänger

Eine natürliche oder juristische Person, Behörde, Einrichtung oder andere Stelle, der personenbezogene Daten offengelegt werden, unabhängig davon, ob es sich bei ihr um einen Dritten handelt oder nicht. Behörden, die im Rahmen eines bestimmten Untersuchungsauftrags nach dem Unionsrecht oder dem Recht der Mitgliedstaaten möglicherweise personenbezogene Daten erhalten, gelten jedoch nicht als Empfänger; die Verarbeitung dieser Daten durch die genannten Behörden

<sup>1</sup> Die hier genannte Definition des Begriffs »anonymisierte Daten« entstammt Erwägungsgrund 26 DSGVO.

erfolgt im Einklang mit den geltenden Datenschutzvorschriften gemäß den Zwecken der Verarbeitung.

**Dritter**

Eine natürliche oder juristische Person, Behörde, Einrichtung oder andere Stelle, außer der betroffenen Person, dem Verantwortlichen, dem Auftragsverarbeiter und den Personen, die unter der unmittelbaren Verantwortung des Verantwortlichen oder des Auftragsverarbeiters befugt sind, die personenbezogenen Daten zu verarbeiten.

## 5 Problemstellung

Eine aktuelle Studie<sup>1</sup> über die Einstellung der Europäer zu den Möglichkeiten von Big Data zeigt, dass 32 Prozent der Befragten behaupten, dass Big Data mehr Vor- als Nachteile bietet, während 51 Prozent der Befragten das Gegenteil denken. Die Studie kommt zu dem Schluss, dass die Zurückhaltung der Befragten weitgehend ausgeräumt werden kann, wenn sie die Vorteile der Nutzung von Big Data für sich selbst und die gesamte Gesellschaft verstehen. Doch was ist Big Data und welche Vorteile können Big-Data-Analysen der Gesellschaft bieten?

Der Begriff Big Data bezeichnet die Erhebung, Analyse und wiederholte Ansammlung großer Mengen von Daten, einschließlich personenbezogener Daten, aus vielfältigen Quellen, die durch Computeralgorithmen und fortschrittliche Datenverarbeitungstechniken automatisiert verarbeitet werden, um bestimmte Korrelationen, Trends und Muster zu ermitteln.<sup>2</sup> Der Grundgedanke von Big Data ist, möglichst viele Informationen aus unterschiedlichen Quellen zusammenzuführen zum Zwecke der Erkennung von Korrelationen und Kausalitäten und der Generierung neuer Informationen.<sup>3</sup> Mit den gewonnenen Informationen können Unternehmen zum Beispiel zielgerichtete Werbung, Marktprognosen sowie die Bekämpfung von Betrug erzielen. Auch für die Bürger bringt die Verwendung von Big Data große Vorteile mit sich – zum Beispiel in den Bereichen Gesundheitsfürsorge (mit Big Data können u. a. Epidemien vorhergesagt werden), Verringerung des Energieverbrauchs (mit Big Data kann analysiert werden, wer wann wie viel Energie benötigt und inwieweit sich die Nachfrage aus Wind und Sonne decken lässt. Unter- und Überversorgungen können vorhergesagt werden), Verbesserung der Verkehrssicherheit (durch die Analyse von Fahrzeugdaten wie Bremsvorgängen sowie Berichte über Verkehrsunfälle können gefährdete Standorte identifiziert werden) sowie die Bekämpfung des Klimawandels (durch die Beobachtung aller relevanten Daten des Schiffsverkehrs in Echtzeit und den Abgleich der Daten mit denen des Zielhafens kann z. B. eine

Echtzeitregulierung der Schiffsgeschwindigkeit erfolgen, der CO<sub>2</sub>-Ausstoß wird erheblich reduziert).

Im Rahmen von Big-Data-Analysen werden in der Regel sowohl personenbezogene Daten (z. B. Nutzerprofile sozialer Netzwerkplattformen, Fahrzeug- und Navigationsgeräteprofile, Klicks im Internet) als auch nicht personenbezogene Daten (z. B. Daten über Luftverschmutzung, Wetterdaten, Straßenauslastung) verwendet.

Wann immer personenbezogene Daten in die Big-Data-Analyse einfließen, haben Unternehmen, die Big-Data-Analysen durchführen, im Rahmen der Analysen die geltenden Datenschutzbestimmungen zu beachten. Dies gilt i. d. R. auch dann, wenn Unternehmen, die Big-Data-Analysen durchführen, gar nicht primär an dem Personenbezug der Ursprungsdaten – also der Daten, welche die Grundlage einer Big-Data-Analyse bilden – selbst interessiert sind, sondern an dem Mehrwert des Analyseergebnisses, das häufig nicht mehr personenbezogen ist bzw. erst wieder Personenbezug besitzt, wenn ein Analysemerkmal einer natürlichen Person zugeordnet wird.

Den europäischen und deutschen Gesetzgeber sowie die Aufsichtsbehörden für den Datenschutz beschäftigt es seit geraumer Zeit, unterschiedliche Interessen auszubalancieren: Einerseits die - meist wirtschaftlichen - Interessen der Big-Data-Analysten, und andererseits die Privatsphäre-Interessen der Personen, deren personenbezogene Daten im Rahmen der Big-Data-Analysen verarbeitet werden.<sup>4</sup>

Die vorliegende Studie hat sich daher zum Ziel gesetzt, den aktuellen Stand der datenschutzrechtlichen Anforderungen an Big Data, die sich seit dem 25. Mai 2018 insbesondere aus der Datenschutz-Grundverordnung und dem neuen Bundesdatenschutzgesetz ergeben, aufzubereiten und Empfehlungen für eine datenschutzkonforme Verwendung von Big Data zu geben.

<sup>1</sup> Vgl. Big Data. A European survey on the opportunities and risks of data analytics. Vodafone Institute für Society and Communications TNS. Januar 2016.

<sup>2</sup> Vgl. [14].

<sup>3</sup> Vgl. [16], Art. 6a) Big Data Rn. 254.

<sup>4</sup> Vgl. [14], [19].



## 6 Abgrenzung personenbezogener, pseudonymisierter und anonymisierter Daten nach DSGVO

Die Verarbeitung personenbezogener Daten ist in der Datenschutz-Grundverordnung (nachfolgend DSGVO) sowie ergänzend in weiteren Datenschutzgesetzen geregelt.<sup>1</sup> Für die Verarbeitung nicht personenbezogener Daten finden die Datenschutzgesetze keine Anwendung. Die Unterscheidung zwischen personenbezogenen Daten und nicht personenbezogenen Daten ist folglich für die Anwendbarkeit der Datenschutzgesetze ausschlaggebend.

Alle Informationen, die sich auf eine identifizierte oder identifizierbare natürliche Person beziehen, sind als personenbezogene Daten zu verstehen, unabhängig von der Form der Information (Sprache, Schrift, Zeichen, Bild oder Ton, digital oder analog).<sup>2</sup> Auch statistische Wahrscheinlichkeitsaussagen, die eine subjektive und/oder objektive Einschätzung zu einer identifizierten oder identifizierbaren Person liefern (z. B. Ausfallwahrscheinlichkeit eines Kredits), weisen einen Personenbezug auf.<sup>3</sup> Dies ist z. B. der Fall, wenn einer Person im Rahmen einer Bonitätsbewertung (z. B. unter Verwendung von Big-Data-Anwendungen) ein sog. Score zugewiesen oder eine Kaufkraftklasse zugeordnet wird.<sup>4</sup> Die Informationen beziehen sich auf identifizierte oder identifizierbare Personen. Eine natürliche Person ist identifizierbar, wenn sie – direkt oder indirekt – mittels Zuordnung zu einer Kennung wie einem Namen, zu einer Kennnummer, zu Standortdaten, zu einer Onlinekennung oder zu einem oder mehreren besonderen Merkmalen bestimmt werden kann, die Ausdruck ihrer physischen, physiologischen, genetischen, wirtschaftlichen, kulturellen oder sozialen Identität sind (vgl. Art. 4 Nr. 1 DSGVO).<sup>5</sup> Nach Erwägungsgrund 26 DSGVO sind bei der Frage, ob eine Person identifizierbar ist, alle Mittel zu berücksichtigen,

»die von dem Verantwortlichen oder einer anderen

Person nach allgemeinem Ermessen wahrscheinlich genutzt werden, um die natürliche Person direkt oder indirekt zu identifizieren, wie beispielsweise das Aussondern. Bei der Feststellung, ob Mittel nach allgemeinem Ermessen wahrscheinlich zur Identifizierung der natürlichen Person genutzt werden, sollten alle objektiven Faktoren, wie die Kosten der Identifizierung und der dafür erforderliche Zeitaufwand, herangezogen werden, wobei die zum Zeitpunkt der Verarbeitung verfügbare Technologie und technologische Entwicklungen zu berücksichtigen sind.«

Für die Identifizierbarkeit einer natürlichen Person genügt es dann, wenn unter Berücksichtigung aller Mittel ohne unverhältnismäßigen Aufwand die Zuordnung hergestellt werden kann. Nicht unumstritten ist, ob der Personenbezug relativ oder absolut zu bestimmen ist.<sup>6</sup> Für die Annahme eines absoluten Personenbezugs genügt es, wenn ein beliebiger Dritter – auch unter Zugrundelegung von Sonderwissen – den Personenbezug herstellen kann. Der relative Personenbezug stellt hingegen darauf ab, dass gerade der Verantwortliche aufgrund seiner Kenntnisse, Mittel und Möglichkeiten den Personenbezug zur betroffenen Person herstellen kann.<sup>7</sup> Die Anwendung des absoluten oder relativen Personenbezugs hat erhebliche Auswirkungen auf die Reichweite des Datenschutzrechts. Die Auslegung des Begriffes ist bspw. bei E-Mail-Adressen, PINs TANs oder IP-Adressen von großer Bedeutung. Der EuGH entschied im Rahmen eines Vorabentscheidungsersuchens des BGH in dem Verfahren Patrick Breyer gegen Bundesrepublik Deutschland, dass bei der Prüfung, ob es sich um ein personenbezogenes Datum handelt, darauf abzustellen sei, ob dem Webseitenbetreiber technische oder rechtliche Mittel zur Verfügung stehen, die es ihm erlauben, die Person

<sup>1</sup> Vgl. Art. 1. Abs. 1 DSGVO, § 1 Abs. 1 BDSG.

<sup>2</sup> Vgl. *Klar/Kühling* in [24], Art. 4 Nr. 1 Rn. 9.

<sup>3</sup> Vgl. [5].

<sup>4</sup> Vgl. *Gola/Klug/Körffner* in [17], § 3 Rn. 3 ff.

<sup>5</sup> Vgl. *Marschall* in [32], § 3 Rn. 8.

<sup>6</sup> Im deutschen Recht begann die Diskussion mit der Auslegung des Begriffs gemäß § 3 Abs. 1 BDSG (alte Fassung) Dazu u. a.: [29]; AG BerlinMitte, 27.03.2007 – 5 C 314/06, ZUM 2008, 83; [13].

<sup>7</sup> Vgl. *Gola/Klug/Körffner* in [17], § 3 Rn. 10.

zu identifizieren. Der EuGH folgte hiermit der Theorie des relativen Personenbezugs.<sup>8</sup> Pseudonymisierte Daten sind folglich für den Verantwortlichen, der die Daten einer bestimmten Person zuordnen kann, personenbezogene Daten. Für Dritte, die nicht über die Zuordnungsregel und/oder über andere Mittel verfügen, sind diese anonym.

## 6.1 Anonymisierung

Unter einer Anonymisierung von personenbezogenen Daten versteht man die Auflösung der Beziehung zwischen den Daten und Angaben und der jeweils betroffenen Person. Gemäß Erwägungsgrund 26 S. 5 der DSGVO gelten die Grundsätze des Datenschutzes nicht für anonyme Informationen, d. h. »(...) für Informationen, die sich nicht auf eine identifizierte oder identifizierbare natürliche Person beziehen, oder personenbezogene Daten, die in einer Weise anonymisiert worden sind, dass die betroffene Person nicht oder nicht mehr identifiziert werden kann.«

Eine absolute Anonymisierung liegt vor, wenn niemand in der Lage ist, den Personenbezug wiederherzustellen. Dies kann bspw. durch die Löschung von Identifikationsmerkmalen in einer Datenbank erfolgen. Da der Verantwortliche das Zusatzwissen Dritter nicht kennen kann, muss er strukturell sicherstellen, dass unabhängig von möglichem Zusatzwissen keine De-Anonymisierung mehr möglich ist.<sup>9</sup> Die absolute Anonymisierung wird gesetzlich bei veröffentlichten amtlichen Statistiken verlangt. (Vgl. § 16 Abs. 1 S. 2 Nr. 4 BStatG)

Kann eine Identifizierung der betroffenen Person nur mit einem unverhältnismäßig großen Aufwand an Zeit, Kosten und Arbeitskraft erfolgen, handelt es sich um eine »faktische Anonymisierung«.<sup>10</sup>

Bei einer absoluten oder faktischen Anonymisierung sind die Datenschutzgesetze nicht anzuwenden. Kann jedoch der Verantwortliche den Personenbezug immer noch wiederher-

stellen, ist die Anonymisierung nicht gelungen und die Daten sind als personenbezogene Daten zu behandeln.<sup>11</sup>

Die Anonymisierung und Pseudonymisierung von Daten wird in der DSGVO als Ausdruck des Grundsatzes der Datenminimierung an mehreren Stellen gefördert. Daten sind zu anonymisieren oder pseudonymisieren, wenn dies nach dem Verwendungszweck möglich ist und in Beziehung zum angestrebten Schutzzweck keinen unverhältnismäßigen Aufwand erfordert (vgl. Art. 5 Abs. 1 und Art. 32 Abs. 1 DSGVO).

stellen, ist die Anonymisierung nicht gelungen und die Daten sind als personenbezogene Daten zu behandeln.<sup>11</sup>

Technische Anonymisierungsverfahren können personenbezogene Eingabedaten in anonyme Ausgabedaten umwandeln (vgl. Kapitel 14), etwa durch Zusammenfassung/Vergrößern von Einzeldaten oder Maskierung von Daten (bspw. durch Einbringen von »künstlichem Rauschen«). Durch die Anonymisierung verlieren die Daten allerdings an Genauigkeit und Aussagekraft. Je nach Anwendungsfall kann es erforderlich sein mehrere Daten zu löschen oder zumindest zu verallgemeinern – je nach Anwendungsfall könnte der anonymisierte Datensatz von »Uta Mustermann, 57 Jahre alt« z. B. »Frau, 50–60 J« lauten. Eine starke Anonymisierung kann dazu führen, dass die Daten nicht mehr aussagekräftig genug und deshalb nicht verwertbar sind. In jedem konkreten Fall ist daher vor der Anonymisierung zu prüfen, welche Datenfelder eines Datensatzes für die weitere Verwendung des anonymisierten Datensatzes besonders wichtig sind und daher im Rahmen des Anonymisierungsprozesses möglichst aussagekräftig erhalten bleiben sollten. Darüber hinaus sind viele Anonymisierungsansätze anfällig für eine De-Anonymisierung mithilfe von zusätzlichen Informationen, die dem Anonymisierenden zum Zeitpunkt der Anonymisierung nicht bekannt waren (oder von ihm möglicherweise nicht beachtet wurden). Demnach sollte das passende Anonymisierungsverfahren unter Berücksichtigung aller Umstände des konkreten Falles – insbesondere der Wahrscheinlichkeit der Re-Identifizierung einzelner Personen – anhand der vorhandenen Daten ausgewählt werden. Dies kann allerdings im Rahmen von Big-Data-Analysen nicht

<sup>8</sup> Vgl. EuGH, 19.10.2016 – C 582/14, NJW 2016, 3579 und BGH, 16.05.2017 – VI ZR 135/13, GRUR 2015, 192.

<sup>9</sup> Vgl. Ziebarth in [35], Rn. 29.

<sup>10</sup> Vgl. Plath/Schreiber in [30], § 3 BDSG Rn. 56, 59.

<sup>11</sup> Vgl. Ziebarth in [35], Rn. 32.

immer sichergestellt werden, zumal Big-Data-Analysen häufig darauf beruhen, laufend neue Daten zum bereits bestehenden Datenbestand hinzuzufügen. Das bereits getroffene Anonymisierungsverfahren kann daher eventuell »mit der Zeit« nutzlos

werden und muss geändert werden.<sup>12</sup> Denn der Verantwortliche ist verpflichtet, die gewählten Maßnahmen zu aktualisieren, wenn sie nicht mehr die Anforderungen der DSGVO erfüllen.<sup>13</sup>

## 6.2 Pseudonymisierung

Die Pseudonymisierung wird in Art. 4 Nr. 5 legal definiert und in Art. 25 und Art. 32 DSGVO als explizites Beispiel für technische und organisatorische Maßnahmen genannt. Die Pseudonymisierung ist

»die Verarbeitung personenbezogener Daten in einer Weise, dass die personenbezogenen Daten ohne Hinzuziehung zusätzlicher Informationen nicht mehr einer spezifischen betroffenen Person zugeordnet werden können, sofern diese zusätzlichen Informationen gesondert aufbewahrt werden und technischen und organisatorischen Maßnahmen unterliegen, die gewährleisten, dass die personenbezogenen Daten nicht einer identifizierten oder identifizierbaren natürlichen Person zugewiesen werden.«

Eine Pseudonymisierung wird eingesetzt, wenn auf die Kenntnis der Identität einer Person verzichtet werden kann und der Personenbezug in bestimmten Fällen wiederher-

gestellt werden können soll. Im Unterschied zu anonymen Daten ist das Bestehen einer Zuordnungsregel charakteristisch für pseudonyme Daten. Die Pseudonymisierung stellt daher eine technische Schutzmaßnahme dar und ist nicht mit der Anonymisierung zu verwechseln.<sup>14</sup> Ein Beispiel für eine Pseudonymisierung ist das Ersetzen von Identifikationsmerkmalen (z. B. eines Namens) durch Kennzeichen (z. B. Personalnummern), um die Bestimmung des Betroffenen auszuschließen oder wesentlich zu erschweren. Der Betroffene bleibt also für denjenigen, der das Pseudonym vergibt, weiterhin feststellbar. Folglich sind die Datenschutzgesetze anwendbar und zu beachten. Art. 89 Abs. 1 DSGVO erfordert, falls möglich, die Pseudonymisierung der Daten bei der privilegierten Verarbeitung von personenbezogenen Daten zu im öffentlichen Interesse liegenden Archivzwecken, zu wissenschaftlichen oder historischen Forschungszwecken und zu statistischen Zwecken.

## 6.3 Datenschutzrechtliche Relevanz von Big-Data-Anwendungen unter Verarbeitung anonymisierter Daten

Daten, die heute durch die aktuellen vorhandenen Techniken anonymisiert werden, können nicht für immer als solche bewertet werden, vgl. Kapitel 17. Es hängt vom technischen Fortschritt und wachsenden Verknüpfungsmöglichkeiten ab, ob einzelne Personen oder Gruppen in einer anonymisierten Big-Data-Datenbank später gegebenenfalls mit Zusatzwissen re-identifiziert werden können.<sup>15</sup> Wenn die Anonymität der Daten nicht zeitunabhängig gesichert werden kann, besteht die Gefahr, dass große Mengen an gespeicherten anonymisierten Daten in der Zukunft durch die Verknüpfung mit

zusätzlichen Informationen oder durch Verwendung neuer Techniken de-anonymisiert werden können. Eine Speicherung von anonymen oder anonymisierten Daten für eine zukünftige Nutzung ist gesetzlich nicht zu beanstanden.<sup>16</sup> Die regelmäßige Prüfung, ob die eingesetzten Anonymisierungstechniken dem aktuellen Stand der Technik entsprechen sowie ob eventuell neue Risiken der Re-Identifizierung vorhanden sind, zeigt sich folglich als unerlässlich.

<sup>12</sup> Vgl. [25] S.429.

<sup>13</sup> Vgl. Art. 24 Abs. 1 S. 2, Art. 32 Abs. 1 Satz 1 lit. d DSGVO.

<sup>14</sup> Vgl. *Klar/Kühling* in [24], Art. 4 Nr. 5 Rn. 2.

<sup>15</sup> Vgl. [25] S. 429.

<sup>16</sup> Nicht zulässig ist lediglich die Vorratsspeicherung von personenbezogenen Daten.

# 7 Regelungen der DSGVO und des neuen Bundesdatenschutzgesetzes im Überblick

Seit 25. Mai 2018 gilt in Europa die Datenschutz-Grundverordnung (nachfolgend DSGVO). Wie bereits oben erläutert, findet die DSGVO nur bei personenbezogenen Daten Anwendung (vgl. Art. 4 Abs. 1 Nr. 1 DSGVO). Im Rahmen von Big-Data-Analysen können Daten verarbeitet und analysiert werden, die nicht personenbezogene Daten sind, z. B. Daten in Bezug auf Umweltverschmutzung oder Klima. In diesen Fällen ist die DSGVO nicht zu beachten. Werden personenbezogene Daten anonymisiert, sodass die Identifizierung einer Person, unter Berücksichtigung aller verfügbaren Mittel, nicht

möglich ist oder nur mit unverhältnismäßigem Aufwand, ist die DSGVO ebenfalls nicht anwendbar. In vielen Fällen ist dennoch eine Identifizierung oder De-Anonymisierung möglich und die Verarbeitung muss folglich den datenschutzrechtlichen Anforderungen genügen. In diesem Abschnitt werden die gesetzlichen Anforderungen erläutert, die einzuhalten sind, wenn die Möglichkeit der Identifizierung einer natürlichen Person oder der De-Anonymisierung der Daten im Rahmen von Big Data besteht.

## 7.1 Erlaubnistatbestände der DSGVO

Eine Verarbeitung personenbezogener Daten ist nur zulässig, wenn mindestens einer der genannten Erlaubnistatbestände in Art. 6 Abs. 1 erfüllt ist. So dürfen personenbezogene Daten nur verarbeitet werden, wenn

- die betroffene Person eingewilligt hat (vgl. Art. 6 Abs. 1 lit. a DSGVO). Die Einwilligung ist nur wirksam, wenn die Anforderungen von Art. 7 DSGVO zusätzlich erfüllt sind. Die Einwilligung muss u. a. frei, explizit und widerrufbar sein.
- die Verarbeitung für die Erfüllung eines Vertrags erfolgt (vgl. Art. 6 Abs. 1 lit. b DSGVO). Hier ist zu berücksichtigen, dass die Verarbeitung personenbezogener Daten nur soweit erfolgt, wie dies objektiv erforderlich ist (z. B. sind im Rahmen eines Kaufvertrages der Name und die Adresse des Käufers notwendig, um die Rechnung oder die Lieferung adressieren zu können). Die Legitimation einer Datenverarbeitung im Rahmen von Big-Data-Analysen aufgrund dieses Erlaubnistatbestandes ist in jedem Einzelfall sorgfältig zu prüfen, denn durch Big-Data-Analyse werden normalerweise die sich aus der Erbringung der Leistung ergebenden Daten analysiert und für andere Zwecke verarbeitet.
- die Verarbeitung zur Erfüllung einer rechtlichen Verpflichtung erforderlich ist (vgl. Art. 6 Abs. 1 lit. c DSGVO). Eine Datenverarbeitung kann aufgrund von Rechtsvorschriften erforderlich sein. Aus Rechtsvorschriften des Handels- und Steuerrechts können sich z. B. umfassende Dokumentations- und Aufbewahrungspflichten ergeben, die erfüllt werden müssen. Diese Rechtsgrundlage ist bei Big-Data-Analysen nicht von Bedeutung.
- die Verarbeitung zum Schutz lebenswichtiger Interessen der betroffenen Person oder anderer Personen erforderlich ist (vgl. Art. 6 Abs. 1 lit. d DSGVO). Hier sind z. B. Fälle gemeint, in denen die Datenverarbeitung für den Schutz des Lebens oder der Gesundheit erforderlich ist, etwa bei einer Naturkatastrophe. Auf dieser Rechtsgrundlage könnten Big-Data-Analysen beruhen, die zum Zweck haben, das Verhalten von natürlichen Personen bei Naturkatastrophen zu erforschen.
- die Verarbeitung zur Wahrung von Aufgaben im öffentlichen Interesse oder zur Ausübung öffentlicher Gewalt, die dem für die Verarbeitung Verantwortlichen übertragen wurde, erfolgt (vgl. Art. 6 Abs. 1 lit. e DSGVO). Die Konkretisierung dieser Rechtsgrundlage erfolgt durch das Unionsrecht oder das Recht der Mitgliedstaaten. In Deutschland ergänzt § 3 BDSG diesen Erlaubnistatbestand. Big-Data-Analysen können auch für die Erfüllung öffentlicher Interessen erforderlich sein. Hier wird die Datenverarbeitung normalerweise durch eine öffentliche Stelle durchgeführt (z. B. »Smart Cities«). Insbesondere regelt Art. 89 Abs. 1 Satz 1 DSGVO:

- »Die Verarbeitung zu im öffentlichen Interesse liegenden Archivzwecken, zu wissenschaftlichen oder historischen Forschungszwecken oder zu statistischen Zwecken unterliegt geeigneten Garantien für die Rechte und Freiheiten der betroffenen Person gemäß dieser Verordnung.«
- die Verarbeitung zur Wahrung berechtigter Interessen des Verantwortlichen oder eines Dritten erforderlich ist (vgl. Art. 6 Abs. 1 lit. f DSGVO). Unter einem »berechtigten Interesse« sind rechtliche, wirtschaftliche und ideelle Interessen zu verstehen. Hier soll eine Abwägung der widerstreitenden Interessen in jedem Einzelfall erfolgen. Dabei dürfen die Interessen oder Grundrechte der betroffenen Person nicht überwiegen. Dieser Erlaubnistatbestand findet auf Behörden keine Anwendung (vgl. Art. 6 Abs. 1 UAbs. 2 DSGVO). Für öffentliche Stellen ist der Erlaubnistatbestand »Wahrung von Aufgaben im öffentlichen Interesse oder zur Ausübung öffentlicher Gewalt« (vgl. Art. 6 Abs. 1 lit. e DSGVO) in diesen Fällen zu prüfen.<sup>1</sup>

Bei der Abwägung der widerstreitenden Interessen sind folgende Punkte zu berücksichtigen:<sup>2</sup>

- Es muss ein berechtigtes Interesse des für die Verarbeitung Verantwortlichen oder eines Dritten bestehen, der ein solches Interesse geltend macht.
- Die Auswirkungen einer solchen Verarbeitung auf die betroffene Person.
- Die Art der zu verarbeitenden Daten und die Form der Verarbeitung.
- Die berechtigten Erwartungen der Betroffenen an die Verarbeitung.
- Das Ungleichgewicht zwischen dem für die Verarbeitung Verantwortlichen und der betroffenen Person

Eine der grundlegenden und vorbereitenden Fragen ist die Existenz eines berechtigten Interesses des Verantwortlichen.

Die Artikel-29-Datenschutzgruppe nennt in ihrer Stellungnahme 06/2014 eine Reihe von Fällen, in denen ein solches Interesse besteht, wie etwa bei der Wahrnehmung des Rechts auf Meinungs- und Informationsfreiheit, auch in den Medien und in der Kunst, Marketing- oder Werbemaßnahmen, Verhinderung von Betrug oder Missbrauch von Dienstleistungen, Sicherheit, Wissenschaft-, Statistik- oder Forschungszwecke. In der Stellungnahme werden auch die Personalisierung von kommerziellen Angeboten und Onlinemarketingaktivitäten erwähnt. Sie warnt jedoch davor, dass selbst wenn ein solches berechtigtes Interesse besteht, es keine ausreichende Grundlage für die Umsetzung von komplexen Kundenprofilen, die einen erheblichen Eingriff in Ihre Privatsphäre darstellen würden, ist. Beim Vorliegen eines berechtigten Interesses sind die Auswirkungen der Verarbeitung auf die Grundrechte und -freiheiten der betroffenen Personen zu prüfen.

Big-Data-Anwender können mehrere legitime Interessen haben wie z. B. die Erstellung von Kundenprofilen, um die Zielgruppe zu erreichen oder die Verhinderung von Betrug und Missbrauch ihrer Dienste. Die Datenverarbeitung ist jedoch nicht erforderlich, wenn es eine andere Art der Verarbeitung zur Erreichung des Zweckes gibt, die die Rechte der betroffenen Person weniger beeinträchtigt. Der Verantwortliche muss seine legitimen Interessen gegen die Interessen der betroffenen Person abwägen und nur, wenn die Interessenabwägung ergibt, dass seine Interessen schwerer wiegen, darf die Verarbeitung stattfinden. Werden mit einer Big-Data-Analyse Persönlichkeitsmerkmale bestimmt, die die betroffene Person selbst nicht aufgedeckt hat, stellt diese Verarbeitung zunächst einen besonders tiefen Eingriff in ihr Grundrecht auf Datenschutz dar.<sup>3</sup> Hier ist mit Sorgfalt zu prüfen, ob die Interessen der betroffenen Person in dem Einzelfall nicht überwiegen. Zu berücksichtigen ist ebenfalls, wie sich die Analyse auf die Lebensqualität der Menschen auswirkt. Die Abwägung kann im konkreten Fall eine komplexe Bewertung sein, die mehrere Personen und Faktoren umfasst. Der Erlaubnistatbestand »Wahrung berechtigter Interessen« ist folglich nicht die »einfachste« zur erfüllenden Rechtsgrundlage für die zulässige Verarbeitung personenbezogener Daten.<sup>4</sup> Je stärker die Auswirkungen auf die betroffenen Personen, desto wichtiger

<sup>1</sup> Vgl. [32], § 4 Rn. 5.

<sup>2</sup> Vgl. [4].

<sup>3</sup> Vgl. *Richter* in [21], Rn. 258.

<sup>4</sup> Vgl. [20], S. 34.

ist die Beachtung entsprechender Schutzmechanismen.<sup>5</sup> Nicht zu vergessen ist ebenfalls, dass die betroffene Person auch über die Datenverarbeitung zu informieren ist, wenn diese

aufgrund »berechtigter Interessen« erfolgt. Die Informationspflicht muss bei Big-Data-Analysen mit personenbezogenen Daten immer erfüllt werden.

## 7.2 Grundsätze des Datenschutzes

Bei der Verarbeitung von personenbezogenen Daten in Big-Data-Analysen sind die allgemeinen Grundsätze des

Datenschutzes zu beachten. Im Folgenden werden sie näher erläutert.

---

### 7.2.1 Verbotprinzip

---

Jede Verarbeitung personenbezogener Daten bedarf einer Rechtsgrundlage, ansonsten ist sie rechtswidrig (vgl. Art. 6 Abs. 1 DSGVO). Die Datenverarbeitung ist nur dann recht-

mäßig, wenn einer der oben genannten Erlaubnistatbestände vorliegt. Im Grundsatz bleibt daher alles verboten, was nicht ausdrücklich erlaubt ist.

---

### 7.2.2 Zweckbindung

---

Ein weiterer Grundsatz im Datenschutz ist der Grundsatz der Zweckbindung. Danach dürfen personenbezogene Daten grundsätzlich nur zu bestimmten Zwecken verarbeitet werden. Dadurch soll verhindert werden, dass personenbezogene Daten zu anderen als den erlaubten Zwecken verwendet werden. Auch Empfänger von personenbezogenen Daten sind an den Zweck gebunden, für den ihnen die Daten übermittelt werden (vgl. Art. 5 Abs. 1 lit. b DSGVO). Der Zweck

des Datenumgangs muss hinreichend bestimmt und bereits vor der Erhebung der Daten festgelegt werden. Dies soll den Informationsfluss für die betroffene Person transparenter machen, damit sie abschätzen kann, wer welche Daten über sie erhält. Das Zweckbindungsprinzip kann eines der Hindernisse für Big-Data-Analysen sein, da der Umfang/Zweck der Verarbeitung nicht immer von vornherein bekannt ist.

---

### 7.2.3 Zweckänderung

---

Eine Zweckänderung ist ausnahmsweise zulässig, wenn die Verarbeitung zu einem neuen Zweck mit dem ursprünglichen Zweck »vereinbar« ist.<sup>6</sup> Die DSGVO besagt, dass die Weiterverarbeitung für im öffentlichen Interesse liegende Archivzwecke, für wissenschaftliche oder historische Forschungszwecke oder für statistische Zwecke nicht als unvereinbar mit dem Erhebungszweck gilt, wenn die Voraussetzungen von Art. 89 Abs. 1 DSGVO erfüllt sind (vgl. Art. 5 Abs. 1 lit. b letzter Satz). Art. 89 Abs. 1 verlangt geeignete Garantien, mittels derer sichergestellt werden kann, dass technische und organisatorische Maßnahmen bestehen, mit denen insbesondere der Grundsatz der Datenminimierung gewährleistet wird. Als

konkretes Beispiel wird die Pseudonymisierung genannt. Bei Big-Data-Analysen werden statistische Methoden eingesetzt. Es ist in dem Einzelfall zu prüfen, ob die Verarbeitung für statistische Zwecke im Sinne der DSGVO erfolgt. Dies wird nur der Fall sein, wenn gemäß Erwägungsgrund 162 DSGVO die Ergebnisse der Big-Data-Analysen, keine personenbezogenen Daten sind, und auch nicht für Maßnahmen gegenüber Einzelnen eingesetzt werden. Bei Big-Data-Analysen wird in der Regel über die personenbezogene Anwendung der Ergebnisse später entschieden, nach dem Bekanntwerden von Korrelationen und Mustern. Hier liegt ebenfalls keine Weiterverarbeitung für statistische Zwecke vor, auch wenn der Big-Data-An-

<sup>5</sup> Vgl. [20], S. 11.

<sup>6</sup> Vgl. *Herbst* in [24], Art. 5 Rn. 24.

wender mit der Datenverarbeitung primär auf das Training der Algorithmen und deren Verbesserung abzielt.<sup>7</sup>

Die Weiterverarbeitung könnte jedoch zulässig sein, wenn der neue Zweck mit dem Erhebungszweck vereinbar ist. In Art. 6 Abs. 4 DSGVO werden einige Kriterien beschrieben, die der Verantwortliche bei der Prüfung der Vereinbarkeit der Zwecke berücksichtigen soll:

- Die Verbindung zwischen den Zwecken, für die die personenbezogenen Daten erhoben wurden, und den Zwecken der beabsichtigten Weiterverarbeitung,
- den Zusammenhang, in dem die personenbezogenen Daten erhoben wurden, insbesondere hinsichtlich des Verhältnisses zwischen den betroffenen Personen und dem Verantwortlichen,
- die Art der personenbezogenen Daten, insbesondere ob besondere Kategorien personenbezogener Daten gemäß Artikel 9 verarbeitet werden oder ob personenbezogene Daten über strafrechtliche Verurteilungen und Straftaten gemäß Artikel 10 verarbeitet werden,

- die möglichen Folgen der beabsichtigten Weiterverarbeitung für die betroffenen Personen,

- das Vorhandensein geeigneter Garantien, wozu Verschlüsselung oder Pseudonymisierung gehören können.

Es handelt sich hier lediglich um Beispiele. Die Berücksichtigung weiterer Kompatibilitätskriterien, insbesondere der in der Stellungnahme der Artikel-29-Datenschutzgruppe zur Zweckbindung beschriebenen Kriterien<sup>8</sup>, kann im Einzelfall erforderlich sein.<sup>9</sup> Die DSGVO besagt, dass statistische Zwecke nicht als mit den ursprünglichen Zwecken unvereinbar angesehen werden dürfen. Jedoch sieht die Verordnung vor, dass der für die Verarbeitung Verantwortliche angemessene

Sicherheitsvorkehrungen treffen muss, um sicherzustellen, dass die Beteiligten nicht identifiziert werden können. Das Gleiche gilt für Big-Data-Analysen zum Zwecke der Wissenschaft oder Forschung.

Für andere »nicht vereinbare« Zwecke dürfen die Daten nur verwendet werden, wenn hierfür eine neue Rechtsgrundlage gegeben ist, d. h. dies eine Rechtsvorschrift erlaubt oder die betroffene Person eingewilligt hat (vgl. Art. 6 Abs. 1 DSGVO).

---

#### 7.2.4 Zeitliche Speicherbegrenzung

---

Die Identifizierung der betroffenen Person bei der Speicherung der Daten darf nur so lange möglich sein, wie es für die Verarbeitungszwecke erforderlich ist (Art. 5 Abs. 1 lit. e DSGVO).<sup>10</sup> Das Merkmal der Erforderlichkeit ist Voraussetzung datenschutzrechtlicher Befugnisnormen (vgl. Art. 6 Abs. 1 lit b, c, e und f DSGVO). Nach dem Grundsatz der Erforderlichkeit soll sich der Umgang mit personenbezogenen Daten auf das beschränken, was für den jeweiligen Zweck des Datenumgangs unbedingt notwendig ist. Die

Verarbeitung von personenbezogenen Daten ist demnach nur erforderlich, wenn der Verantwortliche ohne deren Kenntnis die Dienstleistung nicht, nicht vollständig, nicht rechtmäßig oder nicht in angemessener Zeit erfüllen könnte. Daraus ergibt sich, dass eine Datenverarbeitung auf Verdacht oder Vorrat grundsätzlich unzulässig ist (»Verbot der Vorratsdatenspeicherung«).<sup>11</sup> Nicht mehr erforderliche Daten sind grundsätzlich zu löschen.

<sup>7</sup> Vgl. Richter in [21], Rn. 264–265.

<sup>8</sup> Vgl. [3].

<sup>9</sup> Vgl. Schulz in [16], Art. 6 Rn. 204.

<sup>10</sup> Vgl. Herbst in [24], Art. 5 Rn. 64.

<sup>11</sup> Vgl. Gola/Klug/Körffer in [17], § 13 Rn. 4.

---

### 7.2.5 Datenminimierung

---

Der Grundsatz der Speicherbegrenzung wird durch den Grundsatz der Datenminimierung (vgl. Art. 5 Abs. 1 lit c DSGVO) ergänzt. Dieser Grundsatz bezieht sich unter anderem auf die datenschutzfreundliche Gestaltung von Datenverarbeitungssystemen (Systemdatenschutz). Nach dem Grundsatz der Datenminimierung besteht die Verpflichtung, möglichst Daten ohne Personenbezug zu verwenden (Daten-

vermeidung) bzw. den Personenbezug von Daten auf das Notwendige zu beschränken z. B. weniger Datenfelder in einem Formular (Datensparsamkeit). Generell dürfen nur die Daten verwendet werden, die für den bestimmten Zweck zwingend erforderlich sind. Auch in Big-Data-Szenarien ist der Verantwortliche angehalten, Gestaltungsmöglichkeiten der Verarbeitung unter Verzicht auf personenbezogene Daten zu prüfen.<sup>12</sup>

---

### 7.2.6 Transparenz

---

Die betroffene Person muss, um ihre Rechte wahrnehmen zu können, die nötige Einsicht in die Situation, also »Kenntnis der Sachlage« zu ihren personenbezogenen Daten haben. Sie muss wissen, dass eine Datenverarbeitung stattfindet und überprüfen können, ob das Bild, das durch die Datenverarbeitung über sie erstellt wurde, auch dem entspricht, das sie von sich heraus preisgeben will. Nur so ist sie in der Lage, ihr Recht auf Datenschutz auszuüben und zu schützen. Die Transparenz soll durch datenschutzrechtliche Benachrichtigungs- und Auskunftspflichten des Verantwortlichen gegenüber der

betroffenen Person gewährleistet werden (vgl. Art. 12, Art. 13, Art. 14 und Art. 15 DSGVO). Darüber hinaus wird die Transparenz auch durch die informierte Einwilligung (vgl. Art. 7 DSGVO) hergestellt. In einer Reihe von Fällen des Big-Data-Einsatzes werden den Verantwortlichen jedoch zum Zeitpunkt der Datenerhebung die Zwecke und Kontexte der nachfolgenden Auswertungen noch unbekannt sein,<sup>13</sup> was die in diesem Zeitpunkt erforderliche Transparenz einschränkt.<sup>14</sup> Die Nichteinhaltung des Transparenzprinzips stellt jedoch einen Verstoß gegen Art. 5 Abs. 1 der DSGVO dar.

---

### 7.2.7 Datenrichtigkeit

---

Personenbezogene Daten müssen sachlich richtig und erforderlichenfalls auf dem neuesten Stand sein (vgl. Art. 5 Abs. 1 lit. d DSGVO). »Sachlich richtig« sind die Daten, wenn sie mit der Realität übereinstimmen. Weiterhin müssen die Daten auf dem neuesten Stand sein, wenn es der Zweck der Datenver-

arbeitung erfordert.<sup>15</sup> Erforderlich ist es z. B. bei einer Speicherung von Zutrittsberechtigungen oder sonstigen Berechtigungen der betroffenen Person oder wenn Entscheidungen mit Rechtswirkungen davon abhängen (z. B. alter Negativvermerk bei der Schufa).

---

### 7.2.8 Integrität und Vertraulichkeit

---

Personenbezogene Daten sind in einer Weise zu verarbeiten, die eine angemessene Sicherheit der personenbezogenen Daten gewährleistet. Der Grundsatz referenziert indirekt die technisch-organisatorischen Schutzmaßnahmen (auch TOM genannt), die ein Verantwortlicher zu ergreifen hat, um die Durchführung der Regelungen des Datenschutzrechts – und

somit den Schutz der betroffenen Person vor Eingriffen in ihre Persönlichkeitsrechte – in technisch-organisatorischer Hinsicht zu gewährleisten.

---

<sup>12</sup> Vgl. [31] S. 436.

<sup>13</sup> Vgl. [22] [36].

<sup>14</sup> Vgl. [31] S. 436.

<sup>15</sup> Vgl. *Herbst* in [24], Art. 5 Rn. 61.



### 7.2.9 Rechenschaftspflicht

Der Verantwortliche hat die oben beschriebenen Grundsätze einzuhalten und ist dazu verpflichtet, die Einhaltung dieser Grundsätze auch nachweisen zu können (vgl. Art. 5 Abs. 2 DSGVO). Das Einhalten ist im Sinne der Vorschrift weitwinkender zu verstehen, als das bloße Berücksichtigen bzw. Beachten

der Datenschutz-Grundsätze. Die Rechenschaftspflicht wird u. a. durch Art. 33 und 34 DSGVO ergänzt, durch welche der Verantwortliche in bestimmten Fällen zur Meldung an die Aufsichtsbehörde (inkl. Dokumentation) und Benachrichtigung an die betroffene Person verpflichtet wird.

### 7.2.10 Einwilligung zur Datenverarbeitung im Rahmen von Big Data

Die betroffene Person kann jederzeit in die Verarbeitung ihrer personenbezogenen Daten einwilligen (Art. 6 Abs. 1 DSGVO, Art. 22 Abs. 2 DSGVO). Gemäß Art. 7 i. V. m. Art. 4 Nr. 11 DSGVO muss es sich um eine informierte Einwilligung handeln, die für einen bestimmten Zweck abgegeben wird. Damit die betroffene Person eine informierte Einwilligung erteilen kann, muss sie Sinn und Reichweite der Entscheidung verstehen, was bei umfangreichen und schwer verständlichen Big-Data-Verarbeitungen schwierig zu gewährleisten ist.<sup>16</sup> Darüber hinaus dürfen die Daten gemäß dem Zweckbindungsprinzip, ohne erneute Einwilligung grundsätzlich nur zu diesem vorab festgelegten Zweck verwendet werden (vgl. Art. 5 Abs. 1 lit. b DSGVO). Big Data charakterisiert sich allerdings dadurch, dass große Informationsmengen zusammengeführt werden. Der »neue Analysezweck« für die Weiterverwendung der Daten wird erst nachträglich (also nach der Ermittlung) festgelegt.<sup>17</sup> Hier stellt sich die Frage, ob die betroffene Person in der Praxis eine wirksame Einwilligung abgeben kann, wenn der Zweck nachträglich nach der Zusammenführung der unter-

schiedlichen Daten definiert wird.<sup>18</sup> In der Praxis wird man eine informierte Einwilligung nur bejahen können, wenn die betroffene Person über den aktuellen Stand der Verarbeitung ihrer personenbezogenen Daten informiert wird. Dies könnte durch regelmäßige Mitteilungen über die aktuelle Verarbeitung ihrer personenbezogenen Daten (aktuelle Datenschutznotizen) erreicht werden. Auch sollte mit der Information über eine o. g. beabsichtigte Zweckerweiterung die Einholung einer neuen bzw. aktualisierten Einwilligung inkl. Hinweis auf die Widerrufsmöglichkeit sichergestellt werden. Dies kann durch entsprechende technische Maßnahmen umgesetzt werden.<sup>19</sup>

Falls eine Einwilligung nicht eingeholt werden kann, könnte die Datenverarbeitung gemäß Art. 6 Abs. 1 lit. f »Wahrung berechtigter Interessen« möglich sein. Dies ist nur der Fall, wenn die Verarbeitung für die Erfüllung eines legitimen Zweckes erforderlich ist, und wenn die Interessen der betroffenen Person nicht überwiegen.

## 7.3 Rechte der betroffenen Person

Der Verantwortliche ist verpflichtet, geeignete Maßnahmen zu treffen, um der betroffenen Person alle Informationen, die sich auf die Verarbeitung beziehen, in präziser, transparenter, verständlicher und zugänglicher Form zu übermitteln. Weiter muss der Verantwortliche der betroffenen Person die Ausübung ihrer Rechte erleichtern bzw. sie sicherstellen (vgl. Art. 12 Abs. 1 und 2 DSGVO). Bei der Verwendung von anonymisierten Daten im Big-Data-Kontext können, durch mehrseitige

Verarbeitungsketten mit einer Vielzahl von Akteuren, die nachträgliche Re-Identifizierung zuvor anonymisierter Daten bei Dritten sowie bei der Analyse entstehende Zufallsfunde entstehen.<sup>20</sup> Auch in diesen Fällen ist die Durchsetzung der Rechte der betroffenen Person zu gewährleisten. Die zentralen Betroffenenrechte sind in den Artikeln 12 ff. der DSGVO normiert. Das BDSG regelt ergänzend die Einschränkung der Betroffenenrechte im Kapitel 2 (§ 32 bis 37 BDSG).

<sup>16</sup> Vgl. [9] S. 226.

<sup>17</sup> Vgl. [1], S. 45.

<sup>18</sup> Vgl. [1], S. 45; [8].

<sup>19</sup> Vgl. [19], S. 97.

<sup>20</sup> gl. [15], S. 18 ff.

### 7.3.1 Recht auf Information

Gemäß Art. 13 und 14 DSGVO hat die betroffene Person das Recht, von dem Verantwortlichen über die Details der Datenverarbeitung informiert zu werden. Die DSGVO unterscheidet danach, ob die Daten bei der betroffenen Person (Art. 13) oder nicht bei dieser erhoben werden (Art. 14). Sinn und Zweck dieser Vorschriften ist es, der betroffenen Person die Abschätzung zu ermöglichen, wer was wann und bei welcher Gelegenheit über sie weiß, um selbstbestimmt über die Geltendmachung ihrer Betroffenenrechte entscheiden zu können. Sie soll insbesondere in die Lage versetzt werden, ihren Auskunftsanspruch (Art. 15 DSGVO) oder ihre Korrekturrechte (Art. 16 DSGVO) wahrnehmen zu können.<sup>21</sup>

Folgende Informationen sind u. a. zu erteilen (vgl. Art. 13 DSGVO):

- die Identität des Verantwortlichen
- die Kontaktdaten der/des Datenschutzbeauftragten
- die Verarbeitungszwecke und Rechtsgrundlage
- die Empfänger bzw. Kategorien von Empfängern
- gegebenenfalls Absicht der Datenübermittlung an ein Drittland (vgl. Art. 13 Abs. 1 lit. f)
- die Dauer der Speicherung
- die Rechte der betroffenen Person gem. Art. 15 ff DSGVO
- die Widerrufbarkeit von Einwilligungen
- das Beschwerderecht bei der Aufsichtsbehörde
- ob die Bereitstellung personenbezogener Daten gesetzl./vertragl. vorgeschrieben bzw. für einen Vertragsabschluss erforderlich ist, ob der o. die Betroffene zur Bereitstellung verpflichtet ist sowie mögliche Folgen einer Nichtbereitstellung

- Informationen zu automatisierten Entscheidungsfindungen und Profiling.

Werden personenbezogene Daten nicht bei der betroffenen Person erhoben, bestehen für den Verantwortlichen ebenfalls die oben erwähnten Informationspflichten (vgl. Art. 14 DSGVO). Zusätzlich muss der Verantwortliche die betroffene Person jedoch darüber aufklären, aus welcher Quelle die personenbezogenen Daten stammen und ob es sich dabei um eine öffentlich zugängliche Quelle handelt (vgl. Art. 14 Abs. 2 lit. f DSGVO). Die Informationspflicht erstreckt sich gemäß Art. 13 Abs. 2 lit. f und Art. 14 Abs. 2 lit. g und dem Auskunftsrecht nach Art. 15 Abs. 1 lit. h auch auf »aussagekräftige Informationen über die involvierte Logik sowie die Tragweite und die angestrebten Auswirkungen einer derartigen Verarbeitung für die betroffene Person«.<sup>22</sup> Dieses Recht sollte keine Rechte anderer Personen, etwa »Geschäftsgeheimnisse oder Rechte des geistigen Eigentums und insbesondere das Urheberrecht an Software« beeinträchtigen. Dies darf jedoch nicht dazu führen, dass die betroffene Person jegliche Auskunft verweigert wird.<sup>23</sup>

Allgemeine Aussagen zur Art der gespeicherten Daten sind für das Recht auf Information ausreichend, da die Information nicht die Auskunft ersetzen soll. Nach dem Zweck der Vorschriften soll die Information vor Beginn der Datenerhebung erfüllt werden. Sollen die Daten für einen anderen Zweck weiterverarbeitet werden als den, für den sie erhoben wurden, ist die betroffene Person vorab über den anderen Zweck und weitere zugehörige Sachverhalte u. a. die Speicherdauer und das Bestehen eines Beschwerderechts zu unterrichten (vgl. Art. 13 Abs. 3 DSGVO Art. 14 Abs. 4 in V. m. Abs. 2 DSGVO).

Ausnahmen von der Informationspflicht sind daher eng und im Zweifel zugunsten der betroffenen Person auszulegen. Die DSGVO sieht in Art. 13 Abs. 4 in Verb. m. Abs. 2 DSGVO bzw. Art. 14 Abs. 5 Ausnahmen von der Benachrichtigung der betroffenen Person vor. Eine Pflicht zur Information besteht danach nicht,

<sup>21</sup> Vgl. [17], § 19a Rn. 5.

<sup>22</sup> Vgl. [25] S. 432.

<sup>23</sup> Vgl. Erwägungsgrund 63 DSGVO.

- wenn die betroffene Person über die Information bereits verfügt oder
- wenn Rechtsvorschriften der Union oder der Mitgliedstaaten die Erlangung oder Offenlegung der Daten ausdrücklich regeln oder sie einer Geheimhaltungspflicht unterwerfen (in Deutschland §§ 29, 32, 33 BDSG)
- bei indirekter Erhebung der Daten, wenn die Information sich als unmöglich erweist oder mit einem unverhältnismäßigen Aufwand verbunden ist.

Das BDSG ergänzt in §§ 29, 32, 33 die Einschränkung der Informationspflicht.

### 7.3.2 Auskunftsrecht

Von der Informationspflicht (Art. 13, 14 DSGVO) ist das Recht auf Auskunft

(Art. 15 DSGVO) abzugrenzen. Während die Informationspflichten proaktiv

zu erfüllen sind, ist der Verantwortliche nur auf Antrag dazu verpflichtet, Auskunft zu erteilen. Darüber hinaus ist die Informationspflicht generischer (ohne konkreten Bezug auf die individuelle betroffene Person) als das Auskunftsrecht.

Gemäß Art. 15 DSGVO kann die betroffene Person von dem Verantwortlichen eine Bestätigung darüber verlangen, ob sie betreffende personenbezogene Daten verarbeitet werden, und wenn dies der Fall ist, welche Daten dies genau sind.<sup>24</sup>

Der Verantwortliche ist daher verpflichtet, die Auskunft zu gewähren. Für die Auskunftserteilung fallen für die betroffene Person keine Kosten an –wenn es sich nicht um offenkundig unbegründete oder insb. exzessive Wiederholungsanträge handelt (vgl. Art. 12 Abs. 5, Art. 15 Abs. 3 DSGVO). Die Auskunftserteilung soll je nach Sachverhalt schriftlich, elektronisch oder mündlich erfolgen, möglichst in Form einer Kopie der personenbezogenen Daten (vgl. Art. 12 Abs. 1 und Art. 15 Abs. 3 DSGVO). Gemäß Art. 12 Abs. 3 DSGVO muss die Auskunftserteilung unverzüglich erfolgen, spätestens aber innerhalb eines Monats. In Ausnahmefällen kann die Monatsfrist überschritten werden. Bei häufigen Wiederholungen des Auskunftsrechts ohne nachvollziehbaren Anlass (vgl. Art. 12 Abs. 5 DSGVO) und bei einer Beeinträchtigung der Rechte des Verantwortlichen oder anderer Personen z. B. bei Geschäftsgeheimnissen oder Daten mit Bezug auf andere Personen (vgl. Art. 15 Abs. 4 DSGVO) kann das Recht auf Erhalt einer Kopie ausnahmsweise verweigert werden.

Der Umfang des Anspruchs umfasst alle über die betroffene Person gespeicherten Daten. Dies bezieht sich auch auf den Zweck der Verarbeitung, die geplante Speicherdauer, die

Herkunft der Daten und den Empfänger oder Kategorien von Empfängern der Daten, die Informationen über die Betroffenenrechte, gegebenenfalls das Bestehen einer automatisierten Entscheidungsfindung sowie bei Datenübermittlung in Drittländer über die insoweit gegebenen Garantien z. B. die Verwendung von Standarddatenschutzklauseln der EU-Kommission (vgl. Art. 15 Abs. 2 in Verb. m. Art. 46 DSGVO).

Der deutsche Gesetzgeber hat weitere Eingrenzungen des Auskunftsrechts in §§ 34, 27 Abs. 2, 28 Abs. 2 und 29 Abs. 1 BDSG geregelt. An alle Ausnahmetatbestände sind strenge Voraussetzungen geknüpft.

In Big-Data-Analysen werden unterschiedliche Daten verarbeitet. Personenbezogene Daten werden mit anonymisierten Daten sowie mit sonstigen Daten (z. B. Wetterinfo) zusammengefügt und ausgewertet. Dieser Datenpool kann das »Herauspicken« der personenbezogenen Daten einer Person zur Gewährleistung ihres Rechts auf Auskunft, Berichtigung sowie Löschung erschweren oder sogar unmöglich machen. Big-Data-Anwender müssen dementsprechend bereits bei der Entwicklung sowie bei der Auswahl von Big-Data-Analysen den Datenschutz berücksichtigen. Eine weitere Herausforderung, der Big-Data-Anwender sich stellen müssen, ist die mögliche De-Anonymisierung von Daten durch die eventuelle Verknüpfung von Daten, die anonyme Daten personenbeziehbar machen. In solchen Fällen sind die Rechte der betroffenen Person ebenfalls zu gewährleisten.

Sofern das Auskunftsverlangen der betroffenen Person verweigert wird, ist dies grundsätzlich zu begründen und zu dokumentieren, es sei denn, der Zweck der Auskunftsverweigerung wäre damit gefährdet (vgl. § 34 Abs. 2 BDSG).

<sup>24</sup> Vgl. [6].

---

### 7.3.3 Berichtigung

---

Es kann für die betroffene Person wichtig sein, dass unrichtige Angaben, die über sie gespeichert wurden, berichtigt werden. Für diese Fälle räumt Art. 16 Abs. 1 DSGVO der betroffenen Person das Recht ein, unrichtige personenbezogene Daten berichtigen zu lassen.

Um unrichtige personenbezogene Daten handelt es sich, wenn sie Informationen über Tatsachen enthalten, die mit der Wirklichkeit nicht übereinstimmen oder nur ein unvollständiges Abbild derselben abgeben.<sup>25</sup>

Falls die betroffene Person die Richtigkeit von personenbezogenen Daten bestreitet und sich weder die Richtigkeit noch die Unrichtigkeit feststellen lässt, ist die Verarbeitung der Daten einzuschränken und dürfen diese Daten nicht verwendet werden (vgl. Art. 18 Abs. 1 lit. a DSGVO). Der Verantwortliche ist dann verpflichtet, die Richtigkeit zu prüfen. Solange diese nicht erwiesen ist, wird die Verarbeitung eingeschränkt.

---

### 7.3.4 Löschung

---

Die Löschung von personenbezogenen Daten hat zum einen den Zweck, im Falle von unzulässig gespeicherten Daten den rechtmäßigen Zustand wiederherzustellen. Zum anderen soll den Grundsätzen der Erforderlichkeit und der Datenminimierung Genüge getan werden, indem Daten, die für ihren vorgesehenen Zweck nicht mehr notwendig sind, irreversibel einer weiteren Verwendung entzogen werden. Eine Löschung ist u. a. vorzunehmen, wenn die Verarbeitung der personenbezogenen Daten unzulässig oder die Kenntnis der Daten zur Aufgabenerfüllung nicht mehr erforderlich ist (vgl. Art. 17 Abs. 1 DSGVO). Die betroffene Person hat nach Art. 17 DSGVO (mit bestimmten Ausnahmen) das Recht, die Löschung ihrer Daten zu verlangen. Eine Ausnahme besteht zum Beispiel, soweit die Verarbeitung zur Ausübung des Rechts auf freie Meinungsäußerung und Information erforderlich ist (vgl. Art. 17 Abs. 3 lit. a DSGVO). Ist die Verarbeitung zur Geltendmachung, Ausübung oder Verteidigung von Rechtsansprüchen erforderlich, kann auch darin eine Ausnahme von der Löschungspflicht bestehen (vgl. Art. 17 Abs. 3 lit. e DSGVO).

Eine Verarbeitung von Daten ist dann unzulässig, wenn im Zeitpunkt der Beurteilung keine wirksame Einwilligung der betroffenen Person oder keine entsprechende gesetzliche Ermächtigungsgrundlage (mehr) vorliegt (vgl. Art. 6 Abs. 1 DSGVO). Eine Löschung muss auch dann erfolgen, wenn die gespeicherten Daten für ihren ursprünglichen konkreten Zweck, nicht mehr erforderlich sind. Damit stellt die Löschung

die letzte Stufe der Zweckbindung dar. Der Verantwortliche muss die Löschung bei Vorliegen eines Lösungsgrundes gem. Art. 17 Abs. 1 DSGVO grundsätzlich auch ohne entsprechendes Begehren der betroffenen Person vornehmen.<sup>26</sup> Daher ist durch technische und organisatorische Maßnahmen sicherzustellen, dass die Löschung personenbezogener Daten erfolgt, sobald diese nicht mehr benötigt werden.

Eine Einschränkung der Löschpflicht bei unverhältnismäßig hohem Aufwand ist in § 35 Abs. 1 BDSG vorgesehen. Die Einschränkung gilt jedoch nur für die nicht automatisierte Datenverarbeitung. Somit können Big-Data-Anwender nicht einwenden, dass das Ausfindigmachen von personenbezogenen Daten einer Person besonders aufwendig und die Bedeutung der gespeicherten Informationen sowie die Risiken besonders gering sind. Vielmehr müssen die Prozesse des Verantwortlichen im Lichte des Art. 25 (Datenschutz durch Technikgestaltung und durch datenschutzfreundliche Voreinstellungen) so gestaltet werden, dass die Löschrechte nach Art. 17 effektiv verwirklicht werden.<sup>27</sup> Eine Besonderheit des Lösungsrechts ist das sogenannte »Recht auf Vergessenwerden« (Art. 17 Abs. 2 DSGVO), wenn der Verantwortliche die zu löschenden Daten öffentlich gemacht hat. In diesem Fall muss der Verantwortliche vertretbare Schritte unternehmen, um die Stellen, die diese Daten verarbeiten, zu informieren, dass die betroffene Person von ihnen die Löschung aller Verweise zu diesen Daten oder von Kopien oder Replikationen verlangt hat.<sup>28</sup>

---

<sup>25</sup> Vgl. *Herbst* in [24], § 16 Rn. 8.

<sup>26</sup> Vgl. *Herbst* in [24], Art. 17 Rn. 47.

<sup>27</sup> Vgl. *Nolte/Werkmeister* in [16], Art. 17 Rn. 33.

<sup>28</sup> Vgl. [12].

---

### 7.3.5 Einschränkung der Verarbeitung

---

Die betroffene Person kann in bestimmten Fällen auch die Einschränkung der Verarbeitung verlangen (vgl. Art. 18 DSGVO). Dies wird in Art. 4 Nr. 3 DSGVO als die Markierung gespeicherter Daten definiert mit dem Ziel, ihre künftige Verarbeitung einzuschränken. Die personenbezogenen Daten dürfen dann nicht mehr nach den allgemeinen Erlaubnistatbeständen verarbeitet werden. Das BDSG sieht in § 35 ergänzend zu Art. 17 und 18 DSGVO bestimmte Fälle vor, in denen eine

Einschränkung der Verarbeitung an Stelle der Löschung tritt. Die Einschränkung ist nicht als bloße Kennzeichnung (Einschränkungsvermerk) zu verstehen. Der Verantwortliche muss vielmehr durch die Einrichtung geeigneter Verfahren bzw. durch technische und organisatorische Maßnahmen sicherstellen, dass eine weitere Verarbeitung der markierten Daten nur noch zu den in Art. 18 Abs. 2 DSGVO genannten Zwecken erfolgen kann.

---

### 7.3.6 Recht auf Datenübertragbarkeit

---

Mit der DSGVO wurde zum ersten Mal das Recht auf Datenübertragbarkeit (Art. 20 DSGVO) eingeführt. Dadurch sollen die betroffenen Personen eine bessere Kontrolle über ihre Daten erlangen und die bei einem Anbieter gespeicherten Daten leicht und ohne Behinderung auf einen anderen Anbieter übertragen können. Die betroffene Person kann ihr Recht auf Datenübertragbarkeit geltend machen, wenn die in Art. 20 Abs. 1 DSGVO festgelegten Voraussetzungen vorliegen:<sup>29</sup>

- Die betroffene Person muss die Daten einem Verantwortlichen bereitgestellt haben. Die Form der Bereitstellung, also ob mündlich, schriftlich, elektronisch oder sonstiger Form, ist unerheblich.
- Die Verarbeitung der Daten durch den Verantwortlichen beruht auf einer Einwilligung gemäß Art. 6 Abs. 1 lit. a bzw. Art. 9 Abs. 2 lit. a oder auf einem Vertrag gemäß Art.

6 Abs. 1 lit. b.

- Die Verarbeitung erfolgt mithilfe automatisierter Verfahren.
- Bei der Ausübung des Rechts auf Datenübertragbarkeit dürfen die Rechte und Freiheiten anderer Personen nicht beeinträchtigt werden.

Sind die Voraussetzungen erfüllt, kann die betroffene Person vom Verantwortlichen den Erhalt der Daten verlangen (vgl. Art. 20 Abs. 1 DSGVO). Die Daten müssen von dem Verantwortlichen in einem strukturierten, gängigen und maschinenlesbaren Format gestellt werden. Hierbei soll der aktuelle Stand der Technik beachtet werden. Die betroffene Person kann auch die direkte Übermittlung (Art. 20 Abs. 1 und 2 DSGVO) der Daten an einen Dritten verlangen, soweit dies technisch machbar ist. Der Verantwortliche muss das Recht kostenlos erfüllen.<sup>30</sup>

---

### 7.3.7 Widerspruch

---

Die betroffene Person hat gegenüber Verantwortlichen die Möglichkeit, einer Verarbeitung ihrer personenbezogenen Daten zur Wahrung berechtigter Interessen gemäß Art. 6 Abs. 1 lit. f DSGVO zu widersprechen (vgl. Art. 21 DSGVO). Hierzu muss sie ihre besonderen persönlichen Gründe anführen. Der Widerspruch muss begründet sein. Dies ist dann der Fall, wenn die Interessenabwägung ergibt, dass das schutzwürdige Interesse der betroffenen Person aufgrund ihrer besonderen

persönlichen Situation gegenüber dem Interesse des Verantwortlichen an der Datenverarbeitung überwiegt, zum Beispiel, wenn die betroffene Person bei der Verarbeitung der Daten in ihrer gesellschaftlichen, sozialen, wirtschaftlichen, rechtlichen oder familiären Situation nachteilig betroffen wäre. Dem Widerspruch muss der Verantwortliche nicht nachkommen, wenn er zwingende schutzwürdige Gründe für die Verarbeitung nachweisen kann oder die Verarbeitung mit der

---

<sup>29</sup> Vgl. [11], S. 4.

<sup>30</sup> Vgl. *Kamla* in [30], Art. 20 Rn. 10.

Geltendmachung, Ausübung oder Verteidigung von Rechtsansprüchen begründet werden kann (vgl. Art. 21 Abs. 1 Satz 2 DSGVO). Der Widerspruch ist dann mit entsprechender Begründung abzulehnen. Das Widerspruchsrecht greift jedoch bei den Fällen, in denen die Verarbeitung Zwecken des Profilings dient (vgl. Art. 21 Abs. 1 DSGVO). Zum Beispiel wenn

mit Big Data ein Profil für Geschäftszwecke von Auskunfteien, Bonitätsprüfungsanbietern, Werbeverbänden u. a. erstellt wird und mit der Wahrung berechtigter eigener oder fremder Interessen begründet wird. Hiermit wird ein »Opt-out« für solche Fälle geschaffen.<sup>31</sup>

---

### 7.3.8 Sonstige Rechte der betroffenen Person

---

Neben den in Art. 12 ff. DSGVO genannten Rechten hat die betroffene Person weitere Rechte, die sie gegenüber dem Verantwortlichen geltend machen kann. Diese können sich aus der DSGVO selbst oder aus dem BDSG oder aus Spezialgesetzen ergeben. Sofern der betroffenen Person bspw. durch einen unzulässigen Umgang mit ihren personenbezogenen Daten ein Schaden entsteht, kommen insbesondere Schadensersatzansprüche in Betracht (vgl. Art. 82 DSGVO, § 83 BDSG).

Zudem kann sich jede betroffene Person an die (zuständige) Aufsichtsbehörde wenden, wenn sie der Ansicht ist, bei der Verarbeitung ihrer personenbezogenen Daten durch den Verantwortlichen in ihren Rechten verletzt worden zu sein (vgl. Art. 77 DSGVO). Sofern die Rechte der betroffenen Person durch den Verantwortlichen nicht oder nicht in der vorgeschriebenen Weise gewährt werden, können diese Rechte zudem gerichtlich eingeklagt werden (vgl. 79 DSGVO).

---

### 7.3.9 Profiling

---

Big-Data-Technologien ermöglichen Profilbildungen sowie automatisierte Einzelentscheidungen auf der Grundlage von Algorithmen.<sup>32</sup> Die DSGVO schränkt mit Art. 22 die Befugnisse des Verantwortlichen bei der Anwendung dieser Technologie ein. Der Verantwortliche muss das Recht der betroffenen Person, nicht einer ausschließlich auf automatisierten Verarbeitung beruhenden Entscheidung unterworfen zu werden, die ihr gegenüber rechtliche Folgen oder eine erhebliche Beeinträchtigung entfaltet, beachten (vgl. Art. 21 Abs. 1 DSGVO). Demnach ist bei Big Data neben der Prüfung der Zulässigkeit der Datenverarbeitung (vgl. Art. 6 Abs. 1 DSGVO) die Prüfung der Zulässigkeit der Auswertung (Datenanalyse) erforderlich.<sup>33</sup> Hier ist zu prüfen, ob das Gebot der menschlichen Intervention bei automatisierten Einzelentscheidungen gewährt wurde. Eine automatisierte Entscheidung liegt bspw. gemäß Erwägungsgrund 71 DSGVO vor, wenn eine automatische Ablehnung eines Online-Kreditanspruchs oder Online-Einstellungsverfahrens ohne jegliches menschliche Eingreifen stattfindet. Zu einer derartigen Verarbeitung zählt auch das Profiling. (vgl. Erwägungsgrund 71 Satz. 2 DSGVO). Das »Pro-

filings« wird in der Praxis oft dazu verwendet, Persönlichkeitsprofile zu erstellen und damit Angebote zu personalisieren. Art. 4 Nr. 4 DSGVO definiert Profiling wie folgt:

»jede Art der automatisierten Verarbeitung personenbezogener Daten, die darin besteht, dass diese personenbezogenen Daten verwendet werden, um bestimmte persönliche Aspekte, die sich auf eine natürliche Person beziehen, zu bewerten, insbesondere um Aspekte bezüglich Arbeitsleistung, wirtschaftliche Lage, Gesundheit, persönliche Vorlieben, Interessen, Zuverlässigkeit, Verhalten, Aufenthaltsort oder Ortswechsel dieser natürlichen Person zu analysieren oder vorherzusagen.«

Bei einer auf Profiling gestützten automatisierten Entscheidung ist demnach zu prüfen, ob diese die betroffenen Personen erheblich beeinträchtigt oder rechtliche Folgen mit sich bringt. Eine erhebliche Beeinträchtigung ist zu bejahen, wenn ein Vertragsschluss mit der betroffenen Person aufgrund eines

<sup>31</sup> Vgl. *Forgó* in [7], Art. 21 Rn. 19.

<sup>32</sup> Vgl. *Martini* in [28], Art. 22 Rn. 8.

<sup>33</sup> Vgl. [18], S. 145.

Profiling verweigert wird (Ablehnung eines Arbeitsvertrages z. B.).<sup>34</sup> Auch in den Fällen, in denen eine automatische individuelle Preisdifferenzierung aufgrund des Profiling der Zahlungsbereitschaft des Kunden stattfindet, kann dies im Einzelfall eine erhebliche Beeinträchtigung des Kunden darstellen. Beim sogenannten Geomarketing ist ebenfalls Art. 22 zu beachten. Hier werden Mobilfunkdaten sowie Daten aus dem öffentlichen Wi-Fi analysiert, um Verkehrsströme zu ermitteln und mit soziodemografischen Daten zu verknüpfen.<sup>35</sup> Diese Informationen können dann für Entscheidungen in der Stadtentwicklung herangezogen werden und können ebenfalls Unternehmen bei Umsatzprognosen und bei der Auswahl des Sortiments helfen.<sup>36</sup> Ausnahmsweise können

automatisierte Einzelentscheidungen zulässig sein (vgl. Art. 22 Abs. 2 DSGVO). Dies ist der Fall, wenn die automatische Entscheidung für die Erfüllung eines Vertrages zwischen der betroffenen Person und dem Verantwortlichen erforderlich ist. z. B. im Rahmen eines Smart Homes, wird ein Unternehmen mit der Durchführung von Energieverbrauchsanalysen beauftragt.<sup>37</sup> Weiter kann eine automatisierte Entscheidung zulässig sein, wenn eine Rechtsvorschrift dies anordnet oder wenn eine ausdrückliche Einwilligung vorliegt. Die Einwilligung soll die Voraussetzungen von Art. 7 wie Art. 4 Nr. 11 DSGVO erfüllen. Die betroffene Person soll in Kenntnis der Sachlage, also hinreichend informiert werden und mit Einsichtsfähigkeit ihre Einwilligung in der Datenverarbeitung freiwillig abgeben.<sup>38</sup>

## 7.4 Technischer und organisatorischer Datenschutz

Im Rahmen seiner Verantwortung und Haftung soll der Verantwortliche geeignete und wirksame Datenschutz- und Datensicherheitsmaßnahmen treffen, um zu gewährleisten, dass die Verarbeitungstätigkeiten im Einklang mit der DSGVO stehen. Dies soll gem. Art. 24 Abs. 1 DSGVO insbesondere durch die Umsetzung technischer und organisatorischer Maßnahmen (TOM) geschehen. Technische Maßnahmen haben Auswirkung auf die eigentliche Datenverarbeitung, so z. B. Steuerung von Software- oder Hardwareprozessen, Verschlüsselung, Passwortsicherung aber u. a. auch bauliche Maßnahmen zur Zutrittsbeschränkung zählen dazu. Organisatorische

Maßnahmen beziehen sich insbesondere auf die äußeren Rahmenbedingungen der Verarbeitung, z. B. Vier-Augen-Prinzip, Protokollierungen von Tätigkeiten, Stichprobenroutinen, Mitarbeiterschulungen, Berichtspflichten.<sup>39</sup> Der Übergang zwischen »technisch« und »organisatorisch« ist oft fließend.<sup>40</sup> Im Rahmen von Big-Data-Analysen müssen Big-Data-Anwender u. a. technische und organisatorische Maßnahmen zum Schutz der Datenansammlungen vor Missbrauch durch unbefugte Dritte, Verlust oder zur Sicherstellung der jederzeit wirksamen Ausübung der Rechte der betroffenen Person treffen.

### 7.4.1 Gemeinsame Bedingungen der Art. 24, 25, 32 DSGVO

Art. 24 DSGVO ist eine Generalklausel, die teilweise durch speziellere Normen der DSGVO konkretisiert wird und daher eng mit diesen verknüpft ist. Spezialnormen sind etwa Art. 25 (Privacy by Design and Default), Art. 32 (Datensicherheit) und Art. 35 (Datenschutz-Folgenabschätzung), auf die nachstehend noch eingegangen wird. Die Vorschriften haben das gemeinsame Ziel, der rechtmäßigen Datenverarbeitung, welches

u. a. die Realisierung der Datenschutzgrundsätze des Art. 5 DSGVO durch TOM beinhaltet.

#### Rechtlicher Rahmen zur Auswahl der technischen und organisatorischen Maßnahmen

Das Ziel der technischen und organisatorischen Maßnahmen ist der Schutz der Rechte und Freiheit der betroffenen Person.

<sup>34</sup> Vgl. [18], S. 151.

<sup>35</sup> Vgl. [1], S. 46.

<sup>36</sup> Vgl. [https://www.wigeogis.com/de/intersport\\_standortanalysen\\_fuer\\_einzelhandel\\_wichtiger\\_denn\\_je](https://www.wigeogis.com/de/intersport_standortanalysen_fuer_einzelhandel_wichtiger_denn_je).

<sup>37</sup> Vgl. [18], S. 152 Rn. 620.

<sup>38</sup> Vgl. *Martini* in [28], Art. 22 Rn. 38.

<sup>39</sup> Für vorstehenden Absatz, *Martini* in [28], Art. 24, Rn. 21 f; *Härting* in [24], Art. 24, Rn. 17.

<sup>40</sup> Vgl. *Wolff* in [34], Rn. 825.

Welche TOM jeweils ergriffen werden müssen, ist situationsbedingt und von diversen Faktoren abhängig, z. B. dem Gefährdungspotenzial der Verarbeitung.<sup>41</sup> Je größer die Risiken für die Rechte und Freiheiten natürlicher Personen, desto höher ist der Schutzbedarf. Zur Bewertung sind Art, Umfang, Umstände und die Zwecke der Verarbeitungstätigkeit sowie Eintrittswahrscheinlichkeiten und Schwere der Risiken zu berücksichtigen (vgl. Art. 24 Abs. 1, 25 Abs. 1, 32 Abs. 1 DSGVO). Art. 25 und 32 DSGVO benennen den Stand der Technik und die Implementierungskosten als weitere Faktoren. Gemäß dem Verhältnismäßigkeitsprinzip müssen die Maßnahmen angemessen, erforderlich und verhältnismäßig sein.<sup>42</sup>

### Nachweispflicht

Dass er wirksame TOM ergriffen hat, muss der Verantwortliche nachweisen können (vgl. Art. 24 Abs. 1, Art. 5 Abs. 2 und Erwgr. 74 DSGVO). Ansonsten gilt die Pflicht als nicht erfüllt.<sup>43</sup> Gem. Art. 42 DSGVO genehmigte Zertifizierungsverfahren<sup>44</sup> können als Nachweis dienen (vgl. Art. 24 Abs. 3, 25 Abs. 3, 32 Abs. 3 DSGVO). Nach den Art. 24 Abs. 3 und 32 Abs. 3 DSGVO kann auch die Einhaltung genehmigter Verhaltensregeln<sup>45</sup> für die Nachweiserbringung hilfreich sein.

### Aktualisierungspflicht

Der Verantwortliche ist verpflichtet, die TOM zu überprüfen

und erforderlichenfalls zu aktualisieren (vgl. Art. 24 Abs. 1 S. 2 DSGVO, Art. 32 Abs. 1 Satz 1 lit. d).

### Auftragsverarbeiter

Der Auftragsverarbeiter ist zwar kein direkter Adressat der Art. 24 und 25, der Verantwortliche darf aber nur mit Auftragsverarbeitern zusammenarbeiten, die hinreichend Garantien dafür bieten, dass die TOM so durchgeführt werden, damit die Verarbeitung den Anforderungen der DSGVO genügt (vgl. Art. 28 Abs. 1 DSGVO). Zudem hat der Auftragsverarbeiter u. a. Art. 32 DSGVO zu befolgen.

### Einbindung des Datenschutzbeauftragten

Ist die Bestellung eines Datenschutzbeauftragten verpflichtend, sollte dieser eingebunden werden. Art. 38 Abs. 1 DSGVO verpflichtet Verantwortliche und Auftragsverarbeiter, den Datenschutzbeauftragten »ordnungsgemäß und frühzeitig in alle mit dem Schutz personenbezogener Daten zusammenhängenden Fragen« einzubinden. Entsprechende Beratungs- und Überwachungspflichten des Datenschutzbeauftragten enthält Art. 39 DSGVO. Für Behörden und andere öffentliche Stellen – außer Gerichten – besteht gem. Art. 37 Abs. 1 lit a DSGVO eine grundsätzliche Bestellpflicht, im privatwirtschaftlichen Big-Data-Kontext (vgl. zur Bestellpflicht insb. Art. 37 Abs. 1 lit. b, c DSGVO, § 38 BDSG n. F.).

41 Art. 24, 25, 32, 35 DSGVO enthalten unbestimmte Rechtsbegriffe: »Rechte und Freiheiten« umfassen alle Grundrechte und Grundfreiheiten der Betroffenen. Das »Risiko« bezeichnet das Gefährdungspotential. Die »Schwere« bezieht sich auf immaterielle und materielle Schäden. Die »Eintrittswahrscheinlichkeit« betrifft die Realisierungsgefahr bzw. den statistischen Erwartungswert, mit dem ein Schadensereignis künftig eintreten könnte (ausführlicher in Martini in [28], Art. 24 Rn. 27–30; Wolff in [34], Rn. 827).

42 Vgl. Martini in [28], Art. 24, Rn. 42 f.

43 Vgl. Wolff in [34], Rn. 828.

44 Zertifizierungsverfahren gem. Art. 42, 43 fördern die Einhaltung der DSGVO und können zur Erfüllung verschiedener Nachweispflichten herangezogen werden. Verantwortliche bzw. Auftragnehmer können freiwillig an entsprechenden Audits teilnehmen. Nur gem. Art. 43 Abs. 1–3 akkreditierte Zertifizierungsstellen und die Aufsichtsbehörden sind zu Zertifizierung befugt. Nationale Zertifizierungsstellen erhalten – nach Akkreditierung durch die Deutsche Akkreditierungsstelle (DAKS) – die Befugnis von der für sie zuständigen Aufsichtsbehörde. Die Zertifikate bzw. Prüfsiegel können entweder nationale oder auch grenzüberschreitende Gültigkeit haben. Neben den o. g. Nachweispflichten (gem. Art. 24, 25, 32 DSGVO) können entsprechende Zertifizierungen auch dem Auftragsverarbeiter als Nachweis der Einhaltung der Anforderungen der DSGVO dienen (vgl. Art. 28 Abs. 5 DSGVO). Unter bestimmten Voraussetzungen können die Zertifizierungen geeignete Garantien der Verantwortlichen oder Auftragsverarbeiter für die Drittlandübermittlung personenbezogener Daten darstellen, die den Datenexport ohne besondere Genehmigung der Aufsichtsbehörden erlauben (vgl. Art. 46 Abs. 2 lit f, 42 Abs. 2 DSGVO). Die Zertifikate sind maximal 3 Jahre gültig und können bei Erfüllung der Kriterien verlängert werden (Vgl. Art. 42 Abs. 7 DSGVO).

45 Verbände u. a. Vereinigungen, die Kategorien von Verantwortlichen oder Auftragsverarbeitern vertreten, können Verhaltensregeln ausarbeiten und deren Einhaltung kontrollieren. Die Kodizes, die nationale oder grenzüberschreitende Gültigkeit haben können, werden von der zuständigen Aufsichtsbehörde geprüft und gegebenenfalls genehmigt. Sie sollen insb. den Besonderheiten einzelner Verarbeitungsbereiche und den Bedürfnissen von Kleinst- bis mittleren Unternehmen Rechnung tragen. Art. 40 Abs. 2 DSGVO enthält Beispiele von Regelungsinhalten. Die Einhaltung der Verhaltensregeln kann neben der Erfüllung der o. g. Pflichten (gem. Art. 24 und 32 DSGVO) dem Auftragsverarbeiter auch bei der Erbringung hinreichender Garantien für eine DSGVO-konforme Datenverarbeitung nützen (vgl. Art. 28 Abs. 5 DSGVO). Ebenso kann die Einhaltung entsprechender Kodizes durch Verantwortliche oder Auftragsverarbeiter bei der Beurteilung der Auswirkungen der Verarbeitungsvorgänge und insb. bei der DSFA berücksichtigt werden (vgl. Art. 35 Abs. 8 DSGVO). Auch können genehmigte Verhaltensregeln – unter bestimmten weiteren Voraussetzungen – geeignete Garantien der Verantwortlichen oder Auftragsverarbeiter für die Drittlandübermittlung personenbezogener Daten darstellen, die den Datenexport ohne besondere Genehmigung der Aufsichtsbehörden erlauben (vgl. Art. 46 Abs. 2 lit e und f, 40 Abs. 3 DSGVO).



---

## 7.4.2 Sicherstellung der Rechtskonformität durch technische und organisatorische Maßnahmen – Artikel 24 DSGVO

---

### Allgemeine Sicherstellungs- und Nachweispflicht

Art. 24 ist die einleitende, allgemeine Vorschrift zur Verantwortung und Haftung des Verantwortlichen, die sich über die nachfolgenden Artikel erstreckt.<sup>46</sup> Der Verantwortliche hat durch TOM sicherzustellen und nachzuweisen, dass die Verarbeitung personenbezogener Daten im Einklang mit der DSGVO erfolgt (vgl. Art. 24 Abs. 1 DSGVO). Die Vorschrift verlangt Maßnahmen auch im Vorfeld der eigentlichen Datenverarbeitung.<sup>47</sup>

### Datenschutzvorkehrungen

Sofern es in einem angemessenen Verhältnis zu den Verarbeitungstätigkeiten steht, müssen gem. Art. 24 Abs. 2 DSGVO die TOM geeignete Datenschutzvorkehrungen umfassen. Damit sind alle Maßnahmen gemeint, die geeignet sind, Verstöße gegen die DSGVO durch inhaltliche oder prozedurale Vorgaben zu vermeiden. Insbesondere sind darunter Datenschutzmanagement und Datenschutzkonzepte sowie weitere, umfassende Compliance-Maßnahmen zu verstehen.<sup>48</sup> Im Big-Data-Kontext sind aufgrund der großen Datenbestände u. a. wirksame TOM zur Umsetzung des Rechts auf Datenportabilität relevant (vgl. Art. 20 DSGVO). Auch Archivierungs- und Löschkonzepte sowie entsprechende Tools zur Umsetzung des »Rechts auf Löschung bzw. Vergessenwerden« (vgl. Art. 17 DSGVO). Siehe dazu auch die Erläuterungen unter Unterab-

schnitt 7.3.4 »Löschung«.

Erwägungsgrund 75 DSGVO nennt Beispiele von Verarbeitungsszenarien die zu physischen, materiellen oder immateriellen Schäden führen könnten: Etwa bei Bewertung persönlicher Aspekte, insbesondere, wenn zur Profilerstellung bzw. -nutzung z. B. die Arbeitsleistung, wirtschaftliche Lage, Gesundheit, persönliche Vorlieben oder Interessen, die Zuverlässigkeit, das Verhalten, der Aufenthaltsort oder Ortswechsel analysiert oder prognostiziert werden. Auch bei Verarbeitung einer großen Menge personenbezogener Daten und einer großen Anzahl von betroffenen Personen, bei Verarbeitung besonderer Kategorien personenbezogener Daten (Art. 9 DSGVO) sowie personenbezogener Daten von Schutzbedürftigen (insb. Kinder).

Als mögliche Schadensereignisse nennt Erwgr. 75 z. B. Diskriminierung, Identitätsdiebstahl oder -betrug, finanzielle Verluste, Rufschädigung, Vertraulichkeitsverlust von personenbezogenen Daten, die einem Berufsgeheimnis unterliegen, unbefugte Aufhebung der Pseudonymisierung oder andere erhebliche wirtschaftliche oder gesellschaftliche Nachteile. Zudem, wenn betroffene Personen um ihre Rechte und Freiheiten gebracht oder daran gehindert werden, ihre personenbezogenen Daten zu kontrollieren.

---

## 7.4.3 Art. 25 Datenschutz durch Technikgestaltung und durch datenschutzfreundliche Voreinstellungen

---

Art. 25 konkretisiert die Vorschriften des Art. 24 mit spezifischen Pflichten zum Datenschutz durch Technikgestaltung (Privacy by Design) und datenschutzfreundliche Voreinstellungen (Privacy by Default). Damit soll erreicht werden, dass bereits bei der Entwicklung neuer technischer Produkte und Dienstleistungen datenschutzrechtliche Standards berücksichtigt werden bzw. dass die jeweiligen Voreinstellungen (z. B. von Nutzerkonten) so eingerichtet sind, dass sie die höchste Sicherheit bieten und so wenig Daten wie möglich erhoben und gespeichert werden.

### Datenschutz durch Technikgestaltung (Privacy by Design)

Nach Art. 25 Abs. 1 hat der Verantwortliche sowohl zum Zeitpunkt der Festlegung der Mittel für die Verarbeitung als auch zum Zeitpunkt der eigentlichen Verarbeitung geeignete TOM zu treffen, um dadurch die Datenschutzgrundsätze wirksam umzusetzen und die notwendigen Garantien in die Verarbeitung aufzunehmen, um den Anforderungen der DSGVO zu genügen. Dazu sollte er ein fachbereichs- bzw. abteilungsübergreifendes Datenschutzkonzept entwickeln, das u. a. die Risikoanalyse, Auswahl, Festlegung und Umsetzung der

<sup>46</sup> Vgl. *Piltz* in [16], Art. 24, Rn. 4.

<sup>47</sup> Vgl. *Wolff* in [34], Rn. 821.

<sup>48</sup> Für vorstehenden Absatz: *Martini* in [28], Art. 24, Rn. 40; *Wolff* in [34], Rn. 830.

TOM, Vorgaben zu Dokumentation, Monitoring, Evaluierung und Aktualisierung der TOM beinhaltet.<sup>49</sup> Der Grundsatz der Datenminimierung sowie Pseudonymisierung werden als Beispiele in Art. 25 Abs. 1 erwähnt.

Der Schutz personenbezogener Daten soll bereits in die Programmierung und Konzipierung der Datenverarbeitungsabläufe und der Datenverarbeitungstechnik integriert und bei deren Entwicklung berücksichtigt werden (z. B. durch sog. »Privacy Enhancing Technologies« und entsprechende organisatorische Maßnahmen.)<sup>50</sup> Dies kann auch für den Verantwortlichen (wirtschaftlich) vorteilhaft sein, weil eine nachträgliche Änderung der Ausgestaltung technologischer Systeme mit hohem Aufwand verbunden sein kann.<sup>51</sup>

### Beispiele für geeignete Maßnahmen

Neben den oben erwähnten TOM kommt z. B. die Anonymisierung infrage. Siehe dazu die Erläuterungen unter Kapitel 6 »Abgrenzung personenbezogener, pseudonymisierter und anonymisierter Daten nach DSGVO«. Entsprechend dem Zweckbindungsgrundsatz können die Daten bei Erhebung durch sog. tagging (elektronische Etiketten) einem bestimmten Zweck zugeordnet werden. Gemäß Erwägungsgrund 78 DSGVO sollten Maßnahmen ergriffen werden, durch die und der betroffenen Person ermöglicht wird, die Verarbeitung der personenbezogenen Daten zu überwachen. Auch an TOM zur Erleichterung der Ausübung der Betroffenenrechte ist zu denken. Dazu gehört die technische Umsetzung des Rechts des Betroffenen, seinen Widerspruch (gem. Art. 21 DSGVO) mittels automatisierter Verfahren auszuüben (vgl. Art. 21 Abs. 5 DSGVO).<sup>52</sup>

### Datenschutzfreundliche Voreinstellungen (Privacy by Default)

Der Verantwortliche hat durch TOM sicherzustellen, dass durch Voreinstellung nur personenbezogene Daten, die für

den jeweiligen bestimmten Zweck erforderlich sind, verarbeitet werden (vgl. Art. 25 Abs. 2 S. 1 DSGVO). Dies gilt für die Datenmenge, den Verarbeitungsumfang, die Speicherfrist und Zugänglichkeit der personenbezogenen Daten (vgl. Art. 25 Abs. 2 S. 2 DSGVO). »Voreinstellungen« sind die Standardeinstellungen, die der Diensteanbieter dem Nutzer vorgibt, bevor dieser mit der Nutzung des Systems beginnt.<sup>53</sup> Die Vorschrift soll unterbinden, dass Datenverarbeiter (z. B. Onlinediensteanbieter) durch die Grundeinstellung des Nutzerprofils mehr Daten erheben, als es für die legitimen Nutzungszwecke erforderlich ist. Sie verbietet u. a., eine Einwilligung in darüber hinausgehende Zwecke als Voreinstellung vorzusehen (z. B. Einwilligung in die Weitergabe der Daten für Werbezwecke).<sup>54</sup> Der Verantwortliche wird mit Art. 25 Abs. 2 DSGVO verpflichtet, die Grundsätze der Zweckbindung und Datenminimierung einzuhalten, auch wenn es ggf. gegen seine eigenen Interessen steht.<sup>55</sup> Er hat die Voreinstellungen auf das für den Verarbeitungszweck erforderliche Maß zu begrenzen. Art 25 Abs. 2 DSGVO schließt aber nicht aus, dass Nutzer die Voreinstellungen (z. B. durch Anklicken eines Kästchens o. a. Auswahlmöglichkeiten) verändern können und damit weniger datenschutzfreundliche Einstellungen wählen.<sup>56</sup> Insbesondere ist gem. Art. 25 Abs. 2 S. 3 sicherzustellen, dass personenbezogene Daten durch Voreinstellungen nicht ohne Eingreifen der betroffenen Person einer unbestimmten Zahl von natürlichen Personen zugänglich gemacht werden.

»Eingreifen« meint bewusstes, eigenes Freischalten der Inhalte durch die betroffene Person. Diese in Art. 25 beispielhaft genannte Anforderung soll es den Betroffenen ermöglichen, den Kreis der Empfänger selbst zu steuern und der unbewussten allgemeinen Zugänglichmachung von Inhalten entgegenwirken (z. B. einer unbewussten Einladung der gesamten Social-Media-Öffentlichkeit, Stichwort »Facebook-Party«).<sup>57</sup> Als Voreinstellung ist der kleinstmögliche Empfängerkreis zu wählen.<sup>58</sup>

<sup>49</sup> Vgl. *Nolte/Werkmeister* in [16], Art. 25 Rn. 20.

<sup>50</sup> Vgl. *Martini* in [28], Art. 25, Rn. 10, mit weiteren Verweisen.

<sup>51</sup> Vgl. *Nolte/Werkmeister* in [16], Art. 25 Rn. 2; *Schaar*, Privacy by Design, abrufbar unter: [http://www.bfdi.bund.de/SharedDocs/Publikationen/22Privacy-ByDesign22.pdf?\\_\\_blob=publicationFile](http://www.bfdi.bund.de/SharedDocs/Publikationen/22Privacy-ByDesign22.pdf?__blob=publicationFile), Stand 07.11.18.

<sup>52</sup> Für vorstehenden Absatz: *Martini* in [28], Art. 25, Rn. 29 ff.

<sup>53</sup> Vgl. *Martini* in [28], Art. 25, Rn. 46c.

<sup>54</sup> Vgl. *Martini* in [28], Art. 25, Rn. 45b.

<sup>55</sup> Vgl. *Wolff* in [34], Rn. 832.

<sup>56</sup> Vgl. *Wolff* in [34], Rn. 841.

<sup>57</sup> Vgl. *Martini* in [28], Art. 25, Rn. 52a, Rn. 52d.

<sup>58</sup> Vgl. *Nolte/Werkmeister* in [16] Art. 25 Rn. 31.

#### 7.4.4 Art. 32 – Gewährleistung der Datensicherheit

Artikel 32 Abs. 1 DSGVO konkretisiert den Grundsatz der Integrität und Vertraulichkeit (vgl. Art. 5 Abs. 1 lit. f DSGVO). Die Regelung verpflichtet den Verantwortlichen und den Auftragsverarbeiter. Sie sollen durch geeignete TOM ein dem Risiko der jeweiligen Verarbeitungssituation angemessenes Datensicherheitsniveau gewährleisten.

Verantwortlicher und Auftragsverarbeiter müssen auch sicherstellen, dass die ihnen unterstellten Personen die Verpflichtungen einhalten und personenbezogene Daten nur nach ihren Anweisungen verarbeiten, es sei denn, sie sind nach EU-Recht oder gesetzlichen Vorgaben eines Mitgliedstaates dazu verpflichtet

(vgl. Art. 32 Abs. 4 DSGVO).<sup>59</sup> Infrage kommen z. B. Verpflichtungserklärungen, Belehrungen und technische Maßnahmen zur Zugangs- und Zugriffsbeschränkung.<sup>60</sup>

##### Maßnahmenkatalog

Art. 32 Abs. 1 Hs. 2 DSGVO beinhaltet beispielhafte TOM, die situationsbedingt unter Berücksichtigung der in Hs. 1 genannten Kriterien umgesetzt werden sollten (vgl. Kriterien unter Abschnitt 7.4.1 »Rechtlicher Rahmen zur Auswahl der TOM«). Der Katalog ist nicht abschließend d. h., je nach Verarbeitungssituation können weitere Maßnahmen erforderlich sein:<sup>61</sup>

- die Pseudonymisierung und Verschlüsselung
- die Vertraulichkeit, Integrität, Verfügbarkeit und Belastbarkeit der Systeme und Dienste sind auf Dauer sicherzustellen
- die Verfügbarkeit der personenbezogenen Daten und der Zugang zu ihnen sind bei einem physischen oder technischen Zwischenfall schnell wiederherzustellen
- Verfahren zur regelmäßigen Überprüfung, Bewertung und Evaluierung der Wirksamkeit der TOM

##### Beispiele für weitere geeignete Maßnahmen

Neben den o. g. Maßnahmen gem. Art. 32 Abs. 1 Hs. 2, können z. B. die folgenden umgesetzt werden:<sup>62</sup>

- Zutrittskontrolle, z. B. durch Schließanlagen, Schlüsselverwaltung, Pförtner
- Zugangskontrolle, z. B. durch personalisierte Nutzerkennungen, Passwortrichtlinien
- Zugriffskontrolle, z. B. durch Berechtigungskonzepte, Protokollierungen
- Weitergabe- bzw. Übermittlungskontrolle, z. B. durch Daten- und Verbindungswegverschlüsselung, Authentifizierung, digitale Signatur
- Eingabe- und Transaktionskontrolle, z. B. durch Protokollierungen, Formatbeschränkungen, Nachvollziehbarkeit der Nutzereingaben durch Zeitstempel
- Verfügbarkeitskontrolle, z. B. durch Backup- und Recovery-Verfahren, Notfallmanagement, unterbrechungsfreie Stromversorgung
- Datentrennungs- und Mandantentrennungskontrolle, z. B. durch mandantenfähige Systeme, Instanziierung in Datenbanken, Archivierungskonzept, Richtlinien

Kriterien zur Beurteilung des Schutzniveaus Art. 32 Abs. 2 DSGVO nennt einige Kriterien zur Beurteilung eines angemessenen Schutzniveaus. Danach sind insbesondere folgende Risiken zu berücksichtigen, die mit der Verarbeitung personenbezogener Daten verbunden sind: Vernichtung, Verlust, Veränderung, unbefugte Offenlegung, unbefugter Zugang. Die Vorschrift unterscheidet ausdrücklich nicht, ob die Schadensereignisse unbeabsichtigt oder unrechtmäßig eintreten.

<sup>59</sup> Der Begriff »unterstellte Personen« geht über eigene Mitarbeiter und Angestellte hinaus.

<sup>60</sup> Vgl. *Piltz* in [16], Art. 32, Rn. 49.

<sup>61</sup> Vgl. *Martini* in [28], Art. 32, Rn. 31.

<sup>62</sup> Vgl. *Müller* in [23], E. II. 1. Definitionen zu den erwähnten Kontrollbereichen sind z. B. in Art. 29 der Richtlinie (EU) 2016/680 zu finden.

#### 7.4.5 Standard-Datenschutzmodell (SDM) und Leitlinien der Artikel-29-Datenschutzgruppe zur Orientierung

Die Aufsichtsbehörden haben eine Erprobungsfassung des »Standard-Datenschutzmodells der Datenschutzkonferenz des Bundes und der Länder« veröffentlicht. Es soll u. a. die Verantwortlichen bei der Umsetzung der rechtlichen Vorgaben der DSGVO in technischen und organisatorischen Maßnahmen unterstützen und einen einheitlichen Prüfungsrahmen für die Aufsichtsbehörden schaffen.<sup>63</sup> Ein zugehöriger Maßnahmenkatalog befindet sich noch in der Erarbeitungsphase. Er soll in einzelne Bausteine (auf Ebene der Daten, IT-Systeme, Prozesse) gegliedert werden. Während der Erprobungsphase werden von einzelnen Aufsichtsbehörden Bausteine veröffentlicht, auf Praxistauglichkeit getestet und gegebenenfalls als verbindlich vom AK Technik veröffentlicht.<sup>64</sup> Sieben (noch) nicht verbindliche Bausteine sind derzeit veröffentlicht.<sup>65</sup> Die wesentlichen Elemente des Modells sind:

- die Überführung der Datenschutzerfordernisse in einen Katalog von Gewährleistungszielen, die in der DSGVO verankert sind,
- die Zerlegung der zu betrachtenden Verfahren in drei Komponenten: Daten, IT-Systeme und Prozesse,
- die Klassifizierung der Daten gemäß der Schutzbedarfsabstufungen »normal«, »hoch« und »sehr hoch«,
- der Ableitung des Schutzbedarfs von IT-Systemen und Prozessen anhand des Schutzbedarfs der Daten,
- ein dazu passender Katalog standardisierter Schutzmaßnahmen.

Das SDM unterscheidet insgesamt sieben Gewährleistungsziele: Datenminimierung, Verfügbarkeit, Integrität, Vertraulichkeit, Nichtverkettung<sup>66</sup>, Transparenz<sup>67</sup> und Intervenierbarkeit<sup>68</sup>.

Diese Ziele müssen durch technische und organisatorische Maßnahmen sichergestellt werden. Den Zielen wird jeweils ein Bündel an Maßnahmen zugeordnet, zum Beispiel:

- Die Datenminimierung kann durch die Reduzierung der Angaben, die zu einer Person erfasst werden, die Einschränkung der Zahl der Stellen oder Personen, welchen die Daten zur Verfügung stehen sowie durch die Implementierung automatischer Sperr- und Löschroutinen erreicht werden.
- Die Verfügbarkeit kann durch Backup-Verfahren, das Vorhalten redundanter Hard- und Software oder Maßnahmen zum Schutz gegen Schadsoftware unterstützt werden.
- Zur Wahrung der Integrität personenbezogener Daten tragen Beschränkungen der Schreib- und Änderungsrechte, Prüfsummen und digitale Signaturen sowie dokumentierte Zuweisung von Berechtigungen und Rollen bei.
- Die Vertraulichkeit kann durch die Implementierung eines sicheren Authentifizierungsverfahrens, Verschlüsselung von gespeicherten oder transferierten Daten, sorgfältige Auswahl des Personals und Schutz gegen äußere Bedrohungen, etwa durch Hacker, erreicht werden.
- Die Nichtverkettung kann durch die organisatorische Trennung der Zugriffsmöglichkeiten, ein dazu passendes Berechtigungskonzept, sichere Authentifizierungsverfahren sowie dem Einsatz von zweckspezifischen Pseudonymen, Anonymisierungsdiensten und geregelten Zweckänderungsverfahren erreicht werden.

63 Vgl. [10], S. 22 ff, abrufbar unter: [https://www.bfdi.bund.de/DE/Datenschutz/Themen/Technische\\_Anwendungen/TechnischeAnwendungenArtikel/Standard-Datenschutzmodell.html](https://www.bfdi.bund.de/DE/Datenschutz/Themen/Technische_Anwendungen/TechnischeAnwendungenArtikel/Standard-Datenschutzmodell.html), Stand 10.12.18.

64 Vgl. [10], S. 42.

65 Die Bausteine sind abrufbar unter: <https://www.datenschutz-mv.de/datenschutz/datenschutzmodell/>, Stand 10.12.18.

66 Dieses Ziel bezeichnet die rechtliche Anforderung, dass personenbezogene Daten nicht zusammengeführt werden dürfen. Eine Zusammenführung kann nur erfolgen, wenn diese zur Erfüllung des Zweckes erforderlich ist oder wenn eine Rechtsgrundlage für die Zusammenführung vorliegt (Art. 5 Abs. 1 lit. b. und Art. 6 Abs. 1 DSGVO).

67 Das Ziel der Transparenz bezeichnet die Pflicht des Verantwortlichen die Verarbeitung personenbezogener Daten für die betroffene Person sowie für die Aufsichtsbehörde nachvollziehbar und transparenter zu machen.

68 Dieses Gewährleistungsziel bezeichnet die rechtliche Anforderung, dass die betroffene Person, die ihr zustehenden Rechte jederzeit wirksam ausüben kann.

- Das Ziel der Transparenz kann durch eine sorgfältige Dokumentation der Verarbeitungstätigkeiten sowie die Protokollierung von Zugriffen und Änderungen unterstützt werden.
- Die Intervenierbarkeit (durch die von der Verarbeitung betroffene Person) kann u. a. durch Verfahren zur dokumentierten Bearbeitung von Störungen, Einrichten einer Kontaktstelle für die betroffene Person sowie die Schaffung differenzierter Einwilligung-, Rücknahme- und Widerspruchsmöglichkeiten sichergestellt werden.

Darüber hinaus sind bei der Prüfung der Einhaltung der Gewährleistungsziele gemäß dem Standard-Datenschutzmodell bei einem Datenverarbeitungsvorgang drei Verfahrenskomponenten zu betrachten: (1) personenbezogene Daten, (2) beteiligte IT-Systeme und (3) organisatorische und personelle Prozesse. Maßstab für die Prüfung ist der Schutzbedarf der Daten, der sich an IT-Systeme und Prozesse vererbt. Das SDM unterscheidet drei Schutzbedarfskategorien: sowie die spezifischen Eintrittswahrscheinlichkeiten und Schwere der Risiken zu berücksichtigen.:

#### »Normal«

Geringer kann der Schutzbedarf personenbezogener Daten und deren Verarbeitung nicht eingestuft werden.

#### »Hoch«

In diese Kategorie fallen personenbezogene Daten und deren Verarbeitung, wenn die betroffenen Personen in ihrer gesellschaftlichen Stellung oder in ihren wirtschaftlichen Verhältnissen erheblich beeinträchtigt werden kann.<sup>69</sup>

#### »Sehr hoch«

Hier liegt eine Gefahr für Leib und Leben oder die persönliche Freiheit der betroffenen Person vor.

Die Schutzbedarfseinstufung nach dem SDM ist u. a. von der zuvor durch den Verantwortlichen festgesetzten Risikostufe abhängig. Diese wird unter Berücksichtigung der Art, des Umfangs, der Umstände und Zwecke der Verarbeitungstätigkeit ermittelt.<sup>70</sup> Je höher der Schutzbedarf ist, desto wirksamer müssen die technischen und organisatorischen Maßnahmen sein, mit denen Daten, IT-Systeme und Prozesse geschützt werden. Hierfür ist eine Risikoanalyse unter Berücksichtigung der möglichen Ursachen und Wahrscheinlichkeiten für Datenschutzverletzungen sowie die Evaluierung der Schutzwirkung der umzusetzenden Maßnahmen erforderlich. Eine Orientierungshilfe für die Bestimmung der Risiken einer Verarbeitung bietet das Kurzpapier Nr. 18 Risiko für die Rechte und Freiheiten natürlicher Personen der unabhängigen Datenschutzbehörden des Bundes und der Länder.<sup>71</sup>

Die Artikel-29-Datenschutzgruppe hat ebenfalls diverse Arbeitspapiere veröffentlicht, z. B. zu den Themen »Datenübertragbarkeit« WP 242, »DSFA« WP 248 und »Transparenz« WP 260.<sup>72</sup> Big-Data-Anwender sollten das Standard-Datenschutzmodell sowie die veröffentlichten Arbeitspapiere der Aufsichtsbehörden als Orientierungshilfe für die Umsetzung der DSGVO in Big-Data-Analysen anwenden.

## 7.5 Datenschutz-Folgenabschätzung

Gemäß Art. 35 Abs. 1 ist vor der Verarbeitung personenbezogener Daten, eine Datenschutz-Folgenabschätzung (DSFA) durchzuführen, wenn die Verarbeitung mit einem hohen Risiko für die betroffenen Personen verbunden ist. Dies kann bspw. der Fall sein, wenn neue Technologien der Datenverarbeitung verwendet werden. Weitere Beispiele sind im Gesetz

zu finden. Demzufolge ist eine solche Untersuchung erforderlich, wenn

- eine systematische oder umfassende Bewertung natürlicher Personen mittels Profiling erfolgen soll (vgl. Art. 35 Abs. 3 lit. b DSGVO),

<sup>69</sup> Vgl. [33], S. 75.

<sup>70</sup> Vgl. [10], S. 31 f.

<sup>71</sup> Abrufbar unter: <https://www.datenschutzzentrum.de/artikel/1225-Kurzpapier-Nr.-18-Risiko-fuer-die-Rechte-und-Freiheiten-natuerlicher-Personen.html>, Stand 10.02.2019.

<sup>72</sup> Die Arbeitspapiere sind abrufbar unter: <https://datenschutz.sachsen-anhalt.de/informationen/internationales/datenschutz-grundverordnung/leitlinien-der-artikel-29-datenschutzgruppe/>, Stand 10.12.18.

- eine Verarbeitung sensibler Daten gemäß Art. 9 Abs. 1 oder Art. 10 DSGVO in großem Umfang durchgeführt wird (vgl. Art. 35 Abs. 3 lit. b DSGVO),
- eine systematische großflächige Videoüberwachung im öffentlichen Raum beabsichtigt wird (vgl. Art. 35 Abs. 3 lit. c DSGVO).

Die drei Beispiele sollen nur zur Orientierung dienen.<sup>73</sup> In Erwägungsgrund 21 der DSGVO sind weitere Beispiele genannt. Demnach ist eine DSFA auch dann erforderlich, wenn Verarbeitungsverfahren eingesetzt werden,<sup>74</sup>

- bei denen den betroffenen Personen die Ausübung ihrer Rechte erschwert wird, was bei wenig transparenten Verfahren zu bejahen ist,
- die nach Auffassung der zuständigen Aufsichtsbehörden wahrscheinlich ein hohes Risiko für die Rechte der betroffenen Personen mit sich bringen.

Gemäß der DSGVO müssen die Aufsichtsbehörden Listen mit Verfahren erstellen und veröffentlichen, für die eine Folgenabschätzung erforderlich oder entbehrlich ist (vgl. Art. 35 Abs. 4 bis 6 DSGVO).<sup>75</sup> Die deutsche Aufsichtsbehörde hat für nicht öffentliche Stellen die aktuelle »Liste der Verarbeitungstätigkeiten, für die eine DSFA durchzuführen ist« und für öffentliche Stellen des Bundes die aktuelle »Liste von Verarbeitungsvorgängen gemäß Artikel 35 Abs. 4 DSGVO« auf ihrer Homepage veröffentlicht.<sup>76</sup> Gemäß erstgenannter Liste ist die Durchführung einer DSFA u. a. erforderlich bei einer Zusammenführung von personenbezogenen Daten aus verschiedenen Quellen und Weiterverarbeitung der so zusammengeführten Daten, sofern die Zusammenführung oder Weiterverarbeitung

- in großem Umfang vorgenommen wird,

- für Zwecke erfolgt, für welche nicht alle der zu verarbeitenden Daten direkt bei den betroffenen Personen erhoben wurden,
- die Anwendung von Algorithmen einschließt, die für die betroffenen Personen nicht nachvollziehbar sind, und
- der Erzeugung von Datengrundlagen dient, die dazu genutzt werden können, Entscheidungen zu treffen, die Rechtswirkung gegenüber den betroffenen Personen entfalten, oder diese in ähnlich erheblicher Weise beeinträchtigen können.

Zudem ist gem. o. g. Liste der deutschen Aufsichtsbehörden z. B. eine DSFA durchzuführen, wenn eine Anonymisierung von besonderen Daten nach Art. 9 DSGVO nicht nur in Einzelfällen zum Zwecke der Übermittlung an Dritte erfolgt. Als Beispiele werden Apothekenrechenzentren oder Versicherungen genannt, die umfangreiche besondere personenbezogene Daten anonymisiert und zu anderen Zwecken selbst verarbeiten oder an Dritte weitergeben.

Die Einwilligung der betroffenen Person befreit nicht von der Verpflichtung zur Durchführung einer Folgenabschätzung.<sup>77</sup> Gemäß Art. 39 Abs. 9 DSGVO soll der Verantwortliche gegebenenfalls die betroffenen Personen in die Folgenabschätzung einbeziehen.

Sobald ein Verfahren der Folgenabschätzung unterliegt, ist der Datenschutzbeauftragte zu konsultieren, sofern benannt (vgl. Art. 35 Abs. 2 DSGVO).

In einer DSFA führt der Verantwortliche vorab eine Bewertung der Auswirkungen der vorgesehenen Datenverarbeitung für den Schutz personenbezogener Daten durch (vgl. Art. 35 Abs. 1 S. 1 DSGVO). Diese Untersuchung soll zumindest folgende Punkte umfassen (vgl. Art. 35 Abs. 7 DSGVO):

<sup>73</sup> Vgl. *Marschall* in [32], § 3 Rn. 170.

<sup>74</sup> Vgl. [18], Rn. 39.

<sup>75</sup> Es wurde auch ein Arbeitspapier der Artikel-29-Arbeitsgruppe mit Leitlinien veröffentlicht, wann eine solche Folgenabschätzung notwendig ist. Guidelines on Data Protection Impact Assessment (DPIA) and determining whether processing is »likely to result in a high risk« for the purposes of Regulation 2016/679, wp248rev.01.

<sup>76</sup> Die Liste für den nicht öffentlichen Stellen wurde von der Datenschutzkonferenz (DSK) erstellt. Dies ist ein Gremium der unabhängigen deutschen Datenschutzaufsichtsbehörden des Bundes und der Länder. Die Liste für nicht öffentliche Stellen ist abrufbar unter [https://www.bfdi.bund.de/SharedDocs/Downloads/DE/Datenschutz/Liste\\_VerarbeitungsvorgaengeDSK.pdf?\\_\\_blob=publicationFile&v=3](https://www.bfdi.bund.de/SharedDocs/Downloads/DE/Datenschutz/Liste_VerarbeitungsvorgaengeDSK.pdf?__blob=publicationFile&v=3).

Die Liste der BfDI für die öffentlichen Stellen des Bundes ist abrufbar unter [https://www.bfdi.bund.de/SharedDocs/Downloads/DE/Datenschutz/Liste\\_Verarbeitungsvorgaenge.pdf?\\_\\_blob=publicationFile&v=3](https://www.bfdi.bund.de/SharedDocs/Downloads/DE/Datenschutz/Liste_Verarbeitungsvorgaenge.pdf?__blob=publicationFile&v=3).

<sup>77</sup> Vgl. *Marschall* in [32], § 3 Rn. 173.

- eine detaillierte und umfassende Beschreibung des geplanten Verarbeitungsverfahrens und des verfolgten Zwecks,
- eine Bewertung der Notwendigkeit und Verhältnismäßigkeit des Verfahrens,
- eine Bewertung der Risiken für die Rechte der betroffenen Personen und
- eine detaillierte Beschreibung der Maßnahmen und Sicherheitsvorkehrungen, durch die der Schutz der personenbezogenen Daten sichergestellt wird.

Es ist im Rahmen der DSFA entscheidend, dass das Verfahren umfangreich dokumentiert (vgl. Art. 35 Abs. 7 DSGVO) und in Form eines Berichts öffentlich zugänglich gemacht wird. Die DSGVO fordert nicht explizit die Veröffentlichung des Berichts, dies ist aber aus Transparenzgründen zu empfehlen. In »Zweifelsfällen« sollte die Aufsichtsbehörde konsultiert werden (Art. 36 Abs. 1 DSGVO) und sollten die in Art. 36 genannten Informationen und Unterlagen zur Verfügung gestellt werden. Die Aufsichtsbehörde hat innerhalb einer Regelfrist von 8 Wochen über die Zulässigkeit des Verfahrens zu entscheiden (vgl. Art. 36 Abs. 2 DSGVO).

## 7.6 Verzeichnis von Verarbeitungstätigkeiten

Gemäß Art. 30 Abs. 1 DSGVO muss jeder Verantwortliche oder Auftragsverarbeiter für jede Verarbeitungstätigkeit ein Verzeichnis führen. Unternehmen oder öffentliche Stellen, die personenbezogene Daten durch Big-Data-Analyse verarbeiten, sind in der Regel zur Führung von Verzeichnissen verpflichtet. Die Beschreibung der Verarbeitungstätigkeiten ist ebenfalls aus Transparenz- und Dokumentationsgründen zu empfehlen. In der DSGVO ist der Begriff »Verarbeitung« definiert (vgl. Art. 4 Nr. 2 DSGVO). Demnach ist eine Verarbeitung jeder mit oder ohne Hilfe automatisierter Verfahren ausgeführte Vorgang oder jede solche Vorgangsreihe im Zusammenhang mit personenbezogenen Daten. Dazu zählen das Erheben, Erfassen, die Organisation, das Ordnen, die Speicherung, Anpassung oder Veränderung, das Auslesen, das Abfragen, die Verwendung, die Offenlegung durch Übermittlung, Verbreitung oder eine andere Form der Bereitstellung, der Abgleich oder die Verknüpfung, die Einschränkung, das Löschen oder die Vernichtung.<sup>78</sup>

Der festzulegende Inhalt des Verzeichnisses richtet sich nach Art. 30 Abs. 1 DSGVO. Danach sind folgende Angaben erforderlich:

- Name und Kontaktdaten des Verantwortlichen, des Vertreters des Verantwortlichen sowie der/des Datenschutzbeauftragten,

- die Zwecke der Verarbeitung,
- eine Beschreibung der Kategorien betroffener Personen und der Kategorien personenbezogener Daten,
- die Kategorien von Empfängern, gegenüber denen die Daten offengelegt werden, einschließlich Empfängern in Drittländern,
- eine geplante Datenübermittlung in Drittstaaten sowie die Dokumentierung geeigneter Garantien,
- Regelfristen für die Löschung der Daten,
- eine allgemeine Beschreibung der technischen und organisatorischen Maßnahmen gemäß Art. 32 Abs. 1 DSGVO.

Der Auftragsverarbeiter ist ebenfalls zur Führung eines Verzeichnisses zu allen Kategorien der von ihm im Auftrag durchgeführten Tätigkeiten der Verarbeitung verpflichtet (vgl. Art. 30 Abs. 2 DSGVO).

<sup>78</sup> Art. 30 Abs. 5 DSGVO regelt, welche Verantwortliche zur Führung des Verzeichnisses verpflichtet sind.

## 7.7 Auftragsverarbeitung

Werden Big-Data-Services z. B. aus der Cloud genutzt, müssen die gesetzlichen Anforderungen der Auftragsverarbeitung ebenfalls erfüllt werden. Gerade für kleine Unternehmen bietet die Big-Data-Analyse im Betriebsmodell Cloud Computing die Möglichkeit, von Big Data zu profitieren ohne Investitionen in Personal, Strom, Kühlung oder Rechenzentrumsplatz vorzunehmen.<sup>79</sup>

Eine Auftragsverarbeitung liegt vor, wenn personenbezogene Daten durch einen Auftragsverarbeiter im Auftrag des Verantwortlichen weisungsgebunden verarbeitet werden (vgl. Art. 28 und Art. 29 DSGVO). Der Auftragsverarbeiter ist verpflichtet, bei der Datenverarbeitung ausschließlich auf Weisung des Verantwortlichen zu handeln (vgl. Art. 29 DSGVO). Der Auftraggeber ist weiterhin verantwortlich für die Daten. Die Auftragsverarbeitung wird vom Gesetzgeber privilegiert. Die personenbezogenen Daten können auf Basis eines Auftragsverarbeitungsvertrages an den Auftragsverarbeiter weitergegeben werden. Es bedarf insbesondere keiner Einwilligung der betroffenen Person oder weiterer gesetzlicher Ermächtigungsgrundlagen. Dafür hat der Verantwortliche vor Beginn der Datenverarbeitung mit dem Auftragsverarbeiter einen Vertrag mit den Mindestanforderungen des Art. 28 Abs. 3 Satz 1 und 2 DSGVO abzuschließen, der schriftlich oder in elektronischer Form abgefasst sein kann. Ein wichtiger Bestandteil des Vertrages ist die Darstellung der erforderlichen Maßnahmen zur Sicherheit in der Verarbeitung nach Art. 32 DSGVO. Weiter bedarf der Auftragsverarbeiter für die Beauftragung von Subunternehmen der schriftlichen Genehmigung des Verantwortlichen (vgl. Art. 28 Abs. 2 DSGVO). Der Verantwortliche ist nicht frei bei der Auswahl des Auftragsverarbeiters. Er darf vielmehr nur solche Einrichtungen beauftragen, die hinreichende Garantien dafür bieten, dass die Anforderungen der DSGVO an den Schutz der Rechte der betroffenen Personen

gewährleistet werden und geeignete technische und organisatorische Maßnahmen umgesetzt sind (vgl. Art. 28 Abs. 1 DSGVO). Diese Beurteilung hängt vom Einzelfall ab und muss insbesondere den Schutzbedarf der personenbezogenen Daten berücksichtigen, also deren Wichtigkeit und Sensibilität für die betroffenen Personen. Das Risiko eines ungewollten Datenverlustes oder auch die Möglichkeit unberechtigter Kenntnisnahme der Daten durch Dritte sowie der dadurch mögliche Schaden für die betroffene Person, muss bei der Bewertung des Schutzbedarfs betrachtet werden. Als Beleg solcher Garantien können genehmigte Verhaltensregeln<sup>80</sup> des Auftragsverarbeiters nach Art. 40 DSGVO oder Zertifizierungen nach Art. 42 DSGVO herangezogen werden (vgl. Art. 28 Abs. 5 DSGVO). Eine explizite Verpflichtung des Verantwortlichen zur Durchführung von Folgekontrollen ist in der DSGVO nicht vorhanden. Allerdings ist die regelmäßige Überprüfung der technischen und organisatorischen Maßnahmen während der Zusammenarbeit mit dem Auftragsverarbeiter u. a. auf Grund der Rechenschaftspflicht des Verantwortlichen und seiner Haftung gegenüber den betroffenen Personen dringend zu empfehlen (vgl. Art. 5 Abs. 2 DSGVO). Der Auftragsverarbeiter wird für die Einhaltung des Datenschutzes im Rahmen der Auftragsverarbeitung in die Pflicht genommen (vgl. Art. 82 Abs. 2 Satz 2 DSGVO). Weiter ist er verpflichtet bei der Datenverarbeitung ausschließlich Personen einzusetzen, die sich zur Verschwiegenheit verpflichtet haben.

Die Auftragsverarbeitung ist von dem sogenannten »Joint Controllership« (gemeinsame Verantwortliche gemäß Art. 26 Abs. 1 DSGVO) zu unterscheiden. Wenn mehrere Unternehmen zusammen arbeiten, etwa auf Plattformen oder in Logistikketten und gemeinsam über den Zweck und die Mittel der Verarbeitung entscheiden, liegt gemeinsame Verantwortlichkeit vor.

<sup>79</sup> z. B. Apache Hadoop, Oracle Data Integration Platform Cloud, Big Data mit Amazon Web Services.

<sup>80</sup> Verhaltensregeln nach Art. 40 DSGVO sind ein Instrument der Selbstregulierung. Branchenverbände oder andere Vereinigungen, die von dieser Möglichkeit Gebrauch machen, können damit die teilweise abstrakten Vorgaben der DSGVO für ihren Geschäftsbereich konkretisieren.



## 7.8 Übermittlung an Drittländer oder an internationale Organisationen

Bei der Übermittlung der personenbezogenen Daten an Drittländer oder an internationale Organisationen, sind zusätzlich die in Art. 44 ff. DSGVO festgelegten Bedingungen einzuhalten. Dies ist z. B. der Fall, wenn die Daten der Big-Data-Analyse bei einem ausländischen Cloud-Anbieter gespeichert werden. Hier hat der Big-Data-Anwender zusätzlich zu den allgemeinen Regelungen (Art. 5 ff. DSGVO) zu prüfen, ob eine Feststellung der Angemessenheit des Datenschutzniveaus im Drittland durch Beschluss der Kommission vorhanden ist. Im Falle des Fehlens einer Feststellung der Angemessenheit des Datenschutzniveaus im Drittland, kann die Datenübermittlung zulässig sein, wenn der Verantwortliche oder der Auftragsverarbeiter geeignete Garantien vorsieht, die ein unzureichendes Datenschutzniveau ausgleichen und die Umsetzung der Rechte der betroffenen Personen sicherstellen (vgl. Art. 46 DSGVO). Dazu zählen z. B. EU-Standarddatenschutzklauseln (vgl. Art. 46 Abs. 2 lit c, d). Die folgenden weiteren Beispiele bedürfen der Einzelgenehmigung oder unterliegen Genehmigungsverfahren der jeweils zuständigen Aufsichtsbehörde: Binding Corporate Rules (vgl. Art. 46 Abs. 2 lit b), genehmigte Vertragsklauseln (vgl. Art. 46 Abs. 3 lit a DSGVO) sowie Verhaltensregeln (vgl. Art. 46 Abs. 2 lit e) und Zertifizierungen (vgl. Art. 46 Abs. 2 lit f). Wenn keine der Voraussetzungen gem. den Art. 45, 46 DSGVO vorliegt, ist zu prüfen, ob ein Ausnahmetatbestand gem. Art. 49 DSGVO infrage kommt. z. B. ist die Übermittlung durch die ausdrückliche Einwilligung

der über bestehende Risiken informierten betroffenen Person unter bestimmten weiteren Voraussetzungen möglich (vgl. Art. 49 Abs. 1 UAbs. 1 lit a, UAbs. 2 DSGVO). Da die Einwilligung aber jeweils »für einen bestimmten Fall« eingeholt werden müsste (vgl. Art. 49 Abs. 1 UAbs. 2 S. 1), wird dieser Ausnahmetatbestand auf viele Big-Data-Anwendungen nicht anwendbar sein. Auch die anderen Ausnahmetatbestände des Art. 49 DSGVO sind nur unter engen Voraussetzungen anwendbar.

Gemäß Art. 47 DSGVO bedürfen die verbindlichen internen Datenschutzvorschriften (Binding Corporate Rules) einer Genehmigung durch die zuständige Aufsichtsbehörde. Ist diese Genehmigung jedoch einmal erteilt, bedürfen Datentransfers auf der Grundlage des Binding Corporate Rules keine weitere Genehmigung mehr.

Die sogenannten EU-Standarddatenschutzklauseln wurden durch die Kommission verabschiedet. Diese können durch Verantwortliche und Auftragsverarbeiter ebenfalls verwendet werden und sofern sie unverändert Vertragsbestandteil werden, bieten sie die Garantie für die Einhaltung eines angemessenen Datenschutzniveaus.

## 8 Empfehlungen für Big-Data-Anwender

Die Einhaltung des Datenschutzes in Big-Data-Analysen ist Grundvoraussetzung für die Nutzung von personenbezogenen Daten. Unternehmen, die Big-Data-Technologien einsetzen (Big-Data-Anwender) müssen regelmäßig prüfen, ob in die Big-Data-Analyse personenbezogene Daten einfließen. Die mögliche De-Anonymisierung von Daten aufgrund der Hinzufügung von neuen Daten und die Entstehung neuer Verknüpfungen muss ebenfalls beachtet werden. Das Ergebnis

einer Big-Data-Analyse, das häufig nicht personenbezogen ist, kann ebenfalls wieder Personenbezug bekommen, wenn ein Analysemerkmal einer natürlichen Person zugeordnet wird. Zum Schutz der Privatsphäre des Einzelnen sind die erforderlichen technischen und organisatorischen Maßnahmen zu treffen. Insbesondere sind im Rahmen von Big-Data-Analysen folgende Maßnahmen von großer Relevanz:

### 8.1 Privacy by Design und by Default

Der Datenschutz ist in Big-Data-Szenarien bereits bei der Entwicklung von Big-Data-Lösungen und Anonymisierungsverfahren zu beachten (vgl. Art. 25 DSGVO).<sup>1</sup> Die Anwendung von Big-Data-Technologien soll auf die sieben Grundprinzipien des »Privacy by Design« (Datenschutz durch Technikgestaltung) aufbauen.<sup>2</sup> Diese sind:

1. Proaktiv, nicht reaktiv; als Vorbeugung und nicht als Abhilfe,
2. Datenschutz als Standardeinstellung,
3. Der Datenschutz ist in das Design – also in Programmierung und Konzipierung der Datenverarbeitungsabläufe und -technik eingebettet,
4. Volle Funktionalität – eine Positivsumme, keine Nullsumme,
5. Durchgängige Sicherheit – Schutz während des gesamten Lebenszyklus,
6. Sichtbarkeit und Transparenz – Für Offenheit sorgen,
7. Die Wahrung der Privatsphäre der Nutzer – Für eine nutzerzentrierte Gestaltung sorgen (vgl. Kapitel 11).

Bereits identifizierte Datenschutzrisiken in vergangenen Big-Data-Analysen können durch neue Designlösungen ebenfalls adressiert werden.<sup>3</sup> Die mit auftretenden Risiken verbundenen Kosten und der Zeitaufwand könnten dadurch minimiert werden. Auch die Rechte der betroffenen Person lassen sich in der Praxis bei Big Data nur durchsetzen, wenn die erforderlichen technischen Maßnahmen bereits bei der Entwicklung der Big-Data-Anwendung getroffen wurden.<sup>4</sup> Technische Tools, die die Rechte der betroffenen Person sicherstellen, sollten in der Entwicklungsphase der Big-Data-Anwendung konzipiert werden.

<sup>1</sup> Vgl. [27].

<sup>2</sup> Vgl. [20], S. 15.

<sup>3</sup> Vgl. [19], S. 72.

<sup>4</sup> Vgl. *Nolte/Werkmeister* in [16], Art. 17 Rn. 33.

## 8.2 Regelmäßige Prüfung der Anonymisierungsverfahren sowie sonstigen technischen und organisatorischen Maßnahmen

Verantwortliche sind grundsätzlich dazu verpflichtet zu überprüfen, ob die eingesetzten Verfahren und sonstigen Maßnahmen noch den Datenschutzanforderungen der DSGVO genügen. Erforderlichenfalls müssen Anpassungen vorgenommen werden (vgl. Art. 24 ff. DSGVO). Siehe dazu Abschnitt 7.4 »Technischer und organisatorischer Datenschutz«.

Für Big-Data-Anwender sind u. a. wirksame Anonymisierungsverfahren wichtig. Die DSGVO sieht kein konkretes Anonymisierungsverfahren vor. Maßgeblich für die Anonymisierung ist es, gemäß Art. 4 Nr. 1 DSGVO in Verbindung mit Erwägungsgrund 26, dass keine direkte oder indirekte Identifizierbarkeit der betroffenen Person mehr möglich ist. Entscheidend ist es, ob der Big-Data-Anwender, der das anonymisierte Datenset verarbeitet, »vernünftigerweise« und »nach allgemeinem Ermessen wahrscheinlich« Mittel für eine Re-Identifizierung besitzt.<sup>5</sup> Die Anonymisierung bezieht sich immer auf eine fest umrissene Datenmenge, die anonymisiert werden soll, d. h. Datentypen, Datenumfang und die Gesamtverteilung der Daten sind zum Zeitpunkt der Anonymisierung bereits festgelegt. Wenn die Daten noch nicht vollständig sind, was bei Big Data der Regelfall ist, müssen zumindest Annahmen über das vollständige Datenset gemacht werden. Darüber hinaus ist die

Zusammenführung von zwei anonymisierten Datensets nicht automatisch anonym, da mehr Informationen einen Personenbezug ermöglichen können. Hierfür ist eine risikobasierte Einschätzung einer möglichen Re-Identifizierung im Einzelfall unter Berücksichtigung aller Umstände durchzuführen. Die Robustheit bzw. Schwäche eines Anonymisierungsverfahrens sollte im Einzelfall anhand von drei Kriterien ermittelt werden.<sup>6</sup>

- Besteht immer noch die Möglichkeit eine Einzelperson zu identifizieren?
- Ist es immer noch möglich, Datensätze einer Einzelperson zuzuordnen?
- Können Informationen über eine bestimmte Person abgeleitet werden?

Neben der regelmäßigen Prüfung, ob das Anonymisierungsverfahren »immer noch«<sup>7</sup> wirksam ist und dem Stand der Technik entspricht, empfiehlt sich jährliche Re-Identifizierungstests durchzuführen.<sup>8</sup> Das Testprogramm könnte im Rahmen der Datenschutz-Folgenabschätzung durchgeführt werden.

## 8.3 Mitarbeiterschulung

Sowohl der Verantwortliche als auch seine Mitarbeiter müssen sich der datenschutzrechtlichen Herausforderungen, die mit Big Data verbunden sind, bewusst sein. Der Big-Data-Anwender muss sich diesen Herausforderungen stellen und Schulungen hinsichtlich Maßnahmen zu ihrer Lösung anbieten. Die Praxis zeigt, dass die Schulung der Mitarbeiter eine der wichtigsten Maßnahmen ist, um den Datenschutz im Unternehmen umzusetzen. Alle Beteiligten eines Big-Data-Projekts sollten in der Lage sein, den Datenschutz frühzeitig anzuwenden, insbesondere Anwendungsentwickler, Business Intelligen-

ce Analysts, Datenbankadministratoren und Systemadministratoren. Weiter sollte der Datenschutzbeauftragte einbezogen werden. Er ist derjenige im Unternehmen, der in der Regel mögliche Risiken für die betroffene Person frühzeitig erkennen kann. Neben den rechtlichen Herausforderungen sind auch die ethischen Herausforderungen frühzeitig zu erkennen.<sup>9</sup> Die Verknüpfung von Informationen aus unterschiedlichen Lebenssituationen führt zum »gläsernen Bürger« und ermöglicht über die ursprüngliche Datenbasis hinausgehende Erkenntnisse.<sup>10</sup> Auch die Nutzung von anonymen Daten auf

<sup>5</sup> Vgl. [26] S. 87.

<sup>6</sup> Vgl. [20], S. 13.

<sup>7</sup> Siehe Punkt 2. Abgrenzung personenbezogener, pseudonymisierter und anonymisierter Daten nach DSGVO. Die Problematik der Anwendung der technischen Anonymisierungsverfahren im Rahmen von Big Data.

<sup>8</sup> Vgl. [26] S. 88.

<sup>9</sup> Vgl. [20], S. 15.

<sup>10</sup> Vgl. [31] S. 436.

vorhandenen Personenprofilen kann immer noch Risiken für die betroffene Person darstellen.<sup>11</sup> Hier soll der Big-Data-Anwender vor der Datenverarbeitung durch Big-Data-Techno-

logien seiner ethischen Verantwortung gerecht werden sowie sich der Wirkungen der Verarbeitung auf die betroffenen Personen bewusst sein.

## 8.4 Interne Datenschutzrichtlinie

Big-Data-Anwender sollten in ihren Big-Data-Analysen entsprechendes Datenschutzkonzept entwickeln sowie die erforderlichen technisch-organisatorische Maßnahmen ergreifen, um es umzusetzen.<sup>12</sup> Gemäß Erwägungsgrund 78 soll der Verantwortliche interne Strategien festlegen, um die Einhaltung der DSGVO nachweisen zu können. Interne Datenschutzrichtlinien und ein Datenschutzkonzept können hierfür etabliert

werden. Sie sollten eine Risikoanalyse, Auswahl, Festlegung und Umsetzung konkreter technischer und organisatorischer Maßnahmen umfassen, insbesondere Anonymisierung, Pseudonymisierung, Verschlüsselung und getrennte Speicherung von Daten. Ebenso sollten Maßgaben zur Dokumentation, systematischen Überwachung, Evaluierung sowie Anpassung einzelner Maßnahmen enthalten sein.<sup>13</sup>

## 8.5 Durchführung einer Datenschutz-Folgenabschätzung

Beim Vorliegen eines hohen Risikos ist der Big-Data-Anwender verpflichtet, vor der Datenverarbeitung eine DSFA durchzuführen (vgl. Art. 35 DSGVO). Dafür ist es jedoch erforderlich, dass die Big-Data-Ergebnisse und Risiken für den Big-Data-Anwender vor der Analyse voraussehbar sind, was sich in der Praxis schwierig darstellt. Die Durchführung einer DSFA ist daher auch zu empfehlen, wenn die Risiken der Big-Data-Analyse nicht voraussehbar sind.<sup>14</sup> Die im Rahmen einer DSFA umfassende Beschreibung der Datenverarbeitung und der

organisatorischen und technischen Maßnahmen sowie die Analyse und Bewertung der Risiken für die betroffenen Personen geben dem Big-Data-Anwender einen besseren Überblick über die geplante Verarbeitung. Der Big-Data-Anwender kann hiermit u. a. im Voraus wissen, ob die geplante Datenverarbeitung das geeignete Mittel für die Erfüllung eines Zweckes ist oder ob die Risiken für die betroffene Person, falls voraussehbar, verhältnismäßig sind.

<sup>11</sup> Vgl. [20], S. 16.

<sup>12</sup> Vgl. Marit Hansen, Big Data und Datenschutz abrufbar unter: [https://www.com-magazin.de/praxis/big-data/big-data-treibt-digitalisierung-vor-an-1535480.html?page=5\\_big-data-und-datenschutz](https://www.com-magazin.de/praxis/big-data/big-data-treibt-digitalisierung-vor-an-1535480.html?page=5_big-data-und-datenschutz).

<sup>13</sup> Vgl. Nolte/Werkmeister in [16], Art. 25 Rn. 20.

<sup>14</sup> Dieser Ansicht: [2].

## 8.6 Anforderungen an IT-Systeme

Um das Ziel der rechtmäßigen Datenverarbeitung zu erreichen, sollten die IT-Systeme u. a. folgende Anforderungen erfüllen:

- Es soll nachvollziehbar und nachweisbar sein, welche personenbezogenen Daten zu welchen Zwecken aufgrund welcher erlaubenden Rechtsgrundlage verarbeitet werden.
- Die Zugangs- und Zugriffsrechte auf die Daten müssen nach dem »Need-to-know-Prinzip« (Kenntnis nur bei Bedarf) beschränkt werden.
- Gesetzliche Lösch- und Aufbewahrungspflichten (die sich aus anderen Gesetzen wie z. B. der Abgabenordnung und dem Handelsgesetzbuch ergeben können) müssen eingehalten werden.
- Die Daten müssen lokalisierbar sein d. h., die Speicherorte müssen bekannt sein.
- Die Daten müssen aus allen Speicherorten gelöscht werden können.
- Die Software muss die Anforderungen der datenschutzfreundlichen Technikgestaltung und Voreinstellungen(»Privacy by Design and by Default«) erfüllen.
- Die Datenintegrität muss gewahrt sein.
- Die TOM selbst müssen im Verzeichnis dokumentiert werden.

## 8.7 Genehmigte Verhaltensregeln und Zertifizierung der Big-Data-Anwendungen gemäß DSGVO

Die DSGVO fördert die Selbstregulierung durch genehmigte Verhaltensregeln (vgl. Art. 40 f.) und Zertifizierungsverfahren (vgl. Art. 42 f.).

Durch die Zertifizierung von Big-Data-Anwendungen können Big-Data-Anwender beweisen, dass sie die Anforderungen der DSGVO erfüllen. Die DSGVO fördert in Art. 42 Abs. 2 die Zertifizierungen. Verantwortlicher und Auftragsverarbeiter können zusätzlich zur Einhaltung der DSGVO-Vorschriften datenschutzspezifische Zertifizierungsverfahren, Siegel oder Prüfzeichen, die gemäß der DSGVO genehmigt worden sind, vorsehen. Die Zertifizierung bietet die Möglichkeit, die bestehenden oder demnächst zu implementierenden Big-Data-Anwendungen von unabhängigen Experten überprüfen zu lassen. Dadurch erhält der Big-Data-Anwender eine Rückmeldung zu möglichem Anpassungs- und Verbesserungsbedarf. Das Ergebnis der Evaluierung kann zu Werbezwecken verwendet werden. Dadurch kann das Vertrauen der Nutzer

zur Preisgabe der Daten gewonnen werden. Zu beachten ist jedoch, dass die Zertifizierung nicht die Verantwortung des Verantwortlichen oder Auftragsverarbeiters für die Einhaltung der Datenschutzvorschriften mindert.<sup>15</sup>

Verhaltensregeln nach Art. 40 DSGVO sind ebenfalls ein Instrument der Selbstregulierung. Branchenverbände oder andere Vereinigungen, die Kategorien von Verantwortlichen oder Auftragsverarbeitern vertreten, können damit die teilweise abstrakten Vorgaben der DSGVO für ihren Geschäftsbereich konkretisieren. Die o. g. Verbände und Vereinigungen können Verhaltensregeln ausarbeiten und deren Einhaltung kontrollieren oder die Kontrollbefugnis an eine unabhängige Stelle übertragen. Art. 40 Abs. 2 DSGVO enthält Fallbeispiele, die in entsprechenden Verhaltenskodizes geregelt werden können. Deren Einhaltung dient dem Nachweis der Pflichterfüllung (z. B. bezüglich der TOM) bzw. der Erbringung hinreichender Garantien (z. B. bei der Drittlandübermittlung).<sup>16</sup>

<sup>15</sup> Siehe auch Fußnote 65.

<sup>16</sup> Siehe auch Fußnote 66.

## 9 Zusammenfassung

Um die Risiken für die Rechte und Freiheiten der betroffenen Personen zu minimieren, müssen Unternehmen bei der Entwicklung und beim Einsatz von Big Data sich mit den datenschutzrechtlichen Herausforderungen auseinandersetzen. Es ist unbestritten, dass die Verarbeitung und Analyse von großen Datenmengen relevante Vorteile für die Gesellschaft hat. Mit Big-Data-Technologien lassen sich bspw. Nebenwirkungen von Medikamenten feststellen, die Ausbreitung von Epidemien vorhersagen oder Betrugsfälle aufdecken. Es besteht kein Risiko für die Privatsphäre des Einzelnen, wenn es sich um anonyme Daten handelt wie z. B. Wetterdaten oder Verkehrsdaten.<sup>1</sup> Bei der Verwendung von anonymisierten Daten können jedoch Risiken für die betroffene Person auftreten u. a. wenn eine Re-Identifizierung zu einem späteren Zeitpunkt nicht ausgeschlossen werden kann. Hier ist regelmäßig zu prüfen, ob eine Re-Identifizierung unter Berücksichtigung aller vorhandenen Informationen möglich ist. Die Durchführung einer Risikobewertung im Rahmen einer Datenschutz-Folgenabschätzung sowie die Einhaltung des Datenschutzes bereits

bei der Entwicklung von Big-Data-Anwendungen ist in diesem Fall ebenfalls zu empfehlen. Auch dann, wenn die Informationen zusammengefasst und anonymisiert werden, kann das Ergebnis der Analyse noch Folgen für den Einzelnen haben.<sup>2</sup> Ein Risiko für die Rechte und Freiheiten der betroffenen Person ist immer dann zu bejahen, wenn Big-Data-Analysen durchgeführt werden, um Personenprofile zu erstellen oder zur Vorhersage des Verhaltens von Personen oder Personengruppen. Hier sind die Regelungen der Datenschutz-Grundverordnung zu beachten, insbesondere die Grundsätze der Transparenz und Datenminimierung, das Recht des Einzelnen, nicht einer ausschließlich auf einer automatisierten Verarbeitung beruhenden Entscheidung unterworfen zu werden (vgl. Art. 22 Abs. 1 DSGVO) sowie die Durchführung einer Datenschutzfolgenabschätzung gemäß Art. 35 Abs. 1 DSGVO. Die Einhaltung des Datenschutzes kann zusätzlich zu einer besseren Akzeptanz von Big-Data-Prozessen in der Gesellschaft führen.

---

<sup>1</sup> Vgl. [20], S. 3.

<sup>2</sup> Vgl. [20], S. 3.

## Literatur

- [1] Clemens Appl, Andreas Ekelhart, Natascha Fenz u. a. *Big Data, Innovation und Datenschutz – Studie für eine DSGVO kompatible Vorgangsweise zur Entwicklung einer Big Data Anwendung*. BMVIT. Sep. 2017.
- [2] Art. 29 Datenschutzgruppe. *Guidelines on Data Protection Impact Assessment (DPIA) and determining whether processing is “likely to result in a high risk” for the purposes of Regulation 2016/679*. 17/EN WP 248 rev.01. 2017.
- [3] Art. 29 Datenschutzgruppe. *Opinion 03/2013 on purpose limitation*. WP 203. 2013.
- [4] Art. 29 Datenschutzgruppe. *Stellungnahme 06/2014 zum Begriff des berechtigten Interesses des für die Verarbeitung Verantwortlichen gemäß Art. 7 der Richtlinie 95/46/EG*. WP 217. 2014.
- [5] Art. 29 Datenschutzgruppe. *Stellungnahme 4/2007 zum Begriff „personenbezogene Daten“*. WP 136. 2007.
- [6] Bayerisches Landesamt für Datenschutzaufsicht. *XVI Das Auskunftsrecht der betroffenen Person – Art. 15 DSGVO*. 2017.
- [7] Stefan Brink und Heinrich Amadeus Wolff. *BeckOK Datenschutzrecht*. 25. Edition. Verlag C.H.BECK München, 2018.
- [8] Fred H. Cate und Viktor Mayer-Schönberger. »Notice and consent in a world of Big Data«. In: *Maurer Faculty* (2013).
- [9] Nicolai Culik und Christian Döpke. »Zweckbindungsgrundsatz gegen unkontrollierten Einsatz von Big-Data-Anwendungen«. In: *Zeitschrift für Datenschutz (ZD)* (2017), S. 226–230.
- [10] Datenschutzkonferenz des Bundes und der Länder. *Erprobungsfassung 1.1. des Standard-Datenschutzmodells*. 2018.
- [11] Deutsche Gesellschaft für Medizinische Informatik, Biometrie und Epidemiologie GMDS e.V. *Hinweise zum Recht auf Datenübertragbarkeit gemäß Art. 20 DSGVO*. 2016.
- [12] Die Bundesbeauftragte für den Datenschutz und die Informationsfreiheit (BfDI). *Datenschutz-Grundverordnung Info 6*. 2016.
- [13] Düsseldorfer Kreis. *Beschluss „Datenschutzkonforme Ausgestaltung von Analyseverfahren für Reichweitenmessung bei Internet-Angeboten“*. 2009.
- [14] Europäisches Parlament. *Folgen von Massendaten für die Grundrechte*. P8-TA(2017)0075. 2017.
- [15] FZI Forschungszentrum Informatik. *Smart Data – Smart Privacy? Impulse für eine interdisziplinär rechtlich-technische Evaluation*. 2015.
- [16] Peter Gola. *Datenschutz-Grundverordnung: DS-GVO Kommentar*. Verlag C. H. Beck, 2018.
- [17] Peter Gola und Rudolf Schomerus. *Bundesdatenschutzgesetz Kommentar*. Verlag C.H. Beck, München, 2015.
- [18] Niko Härting. *Datenschutz-Grundverordnung – Das neue Datenschutzrecht in der betrieblichen Praxis*. Verlag Dr. Otto Schmidt, Köln, 2016.
- [19] Information Commissioner’s Office ICO. *Big Data, artificial intelligence, machine learning and data protection*. 2018.
- [20] International Working Group on Data Protection in Telecommunications. *Arbeitspapier zu Big Data und Datenschutz Bedrohung der Grundsätze des Datenschutzes in Zeiten von Big-Data-Analysen*. 675.48.15. 2014.
- [21] Silke Jandt und Roland Steidle, Hrsg. *Datenschutz im Internet – Rechtshandbuch zu DSGVO und BDSG*. 2018.

- [22] Peter Katko und Ayda Babaei-Beigi. »Accountability statt Einwilligung? Führt Big Data zum Paradigmenwechsel im Datenschutz?« In: *Multimedia und Recht (MMR): Zeitschrift für Informations-, Telekommunikations- und Medienrecht*, (2014), 360 ff.
- [23] Ansgar Koreng und Matthias Lachenmann. *Formularhanbuch Datenschutzrecht*. Verlag C.H. Beck, 2018.
- [24] Jürgen Kühling und Benedikt Buchner. *Datenschutz-Grundverordnung Kommentar*. C.H.Beck, 2017.
- [25] Ninja Marnau. »Anonymisierung, Pseudonymisierung und Transparenz für Big Data«. In: *Datenschutz und Datensicherheit (DuD)* 40.7 (Juli 2016), S. 428–433. ISSN: 1862-2607.
- [26] Ninja Marnau, Pascal Berrang und Mathias Humbert. »Anonymisierungsverfahren für genetische Daten«. In: *Datenschutz und Datensicherheit (DuD)* 42.2 (Feb. 2018), S. 83–88. ISSN: 1862-2607.
- [27] Jan-Peter Ohrtmann und Sebastian Schwiering. »Big Data und Datenschutz – Rechtliche Herausforderungen und Lösungsansätze«. In: *Neue Juristische Wochenschrift (NJW)* (2014), S. 2984–2989.
- [28] Boris Paal und Daniel Pauly. *Datenschutz-Grundverordnung / Bundesdatenschutzgesetz: DSGVO – BDSG*. Verlag C. H. Beck, München, 2018.
- [29] Ingrid Pahlen-Brandt. »Datenschutz braucht scharfe Instrumente«. In: *Datenschutz und Datensicherheit (DuD)* 32.1 (2008), S. 34–40.
- [30] Kai-Uwe Plath. *BDSG/DSGVO: Kommentar zum BDSG und zur DSGVO sowie den Datenschutzbestimmungen des TMG und TKG*. Verlag Dr. Otto Schmidt, Köln, 2016.
- [31] Oliver Raabe und Manuela Wagner. »Verantwortlicher Einsatz von Big Data«. In: *Datenschutz und Datensicherheit (DuD)* 40.7 (Juli 2016), S. 434–439. ISSN: 1862-2607.
- [32] Alexander Roßnagel. *Europäische Datenschutz-Grundverordnung, Vorrang des Unionsrechts – Anwendbarkeit des nationalen Rechts*. Nomos Verlagsgesellschaft, Baden-Baden, 2017.
- [33] Markus Schaeffter. *Verfahrensverzeichnis 2.0 – Datenschutzdokumentation konform zur EU-Datenschutzgrundverordnung gestalten*. 2016.
- [34] Peter Schantz und Henrich Amadeus Wolf. *Das neue Datenschutzrecht*. C. H. Beck, 2017.
- [35] Gernot Sydow. *Europäische Datenschutzgrundverordnung Handkommentar*. MANZ Verlag Wien, 2018.
- [36] Thilo Weichert. »Big Data und Datenschutz Chancen und Risiken einer neuen Form der Datenanalyse«. In: *Zeitschrift für Datenschutz (ZD)* (2013), 251 ff.



### III TECHNISCHE ANSÄTZE ZUM SCHUTZ DER PRIVATSPHÄRE



## 10 Schutzziele bei Big Data

In diesem Kapitel werden die verschiedenen Sicherheitsprobleme und Anforderungen behandelt, die beim Umgang mit Big Data auftreten können. Unsere Darstellung folgt zwei Perspektiven. Die erste Perspektive beginnt mit den Grundmotiven für die Absicherung von Daten. Sie berücksichtigt die Datenschutzerfordernisse beim Umgang mit personen-

bezogenen Daten und den Schutz des aus den verarbeiteten Daten gewonnenen Wissens. Die zweite Perspektive folgt den Verarbeitungsschritten der Daten in einem Big-Data-System. Dabei geht es um sichere Übertragungskanäle, sichere Datenspeicherung sowie die Sicherheit bei der Datenverarbeitung.

### 10.1 Gründe für die Absicherung von Daten und Anforderungen in Big-Data-Szenarien

Die Absicherung eines Big-Data-Systems lässt sich wie folgt motivieren: zum einen mit dem Schutz der Personen, zu denen Daten im System gespeichert sind, und zum anderen mit dem Schutz von Wissen – seien es die Rohdaten, die Verarbeitungsalgorithmen oder bereits gewonnene Erkenntnisse als Vorsprung vor Wettbewerbern oder etwa ermittlung-

taktisches oder nachrichtendienstliches Wissen in Sicherheitsbehörden. Oft sind in der Praxis sogar beide Aspekte relevant. Nachfolgend werden beide Motive ausgeführt, aber im Rest der Studie wird der Datenschutz als treibendes Motiv betrachtet.

---

#### 10.1.1 Schutz der Privatsphäre

---

Der Datenschutz zielt auf den Schutz der Privatsphäre von Personen ab, und er ist ein sehr wichtiger Grund für die Sicherung von Big Data. Er ist das Kernthema dieses Berichts, und in den folgenden Kapiteln gehen wir auf einige Technologien zum Schutz der Privatsphäre ein.

Der Schutz der Privatsphäre ist ein grundlegender Aspekt in der heutigen Gesellschaft. Sie ist in vielen normativen Beschlüssen und Rechtsakten wie der Allgemeinen Erklärung der Menschenrechte [57], der Europäischen Menschenrechtskonvention [32] und der Charta der Grundrechte der Europäischen Union verankert und wird z. B. für das »Privat- und Familienleben« allgemein akzeptiert und beantragt. Die seit dem 25. Mai 2018 geltende europäische Datenschutzgrundverordnung (DSGVO), welche in Teil II bereits ausführlich inhaltlich behandelt wurde, ist ein neuer Meilenstein in der Geschichte des Datenschutzes und hat das Potenzial, den Datenschutz insbesondere in der digitalen Welt wirksam durchzusetzen.

In den letzten zwei Jahrzehnten hat die Verbreitung und Nutzung von Daten rasant zugenommen, aber die Rechte auf Privatsphäre und Datenschutz wurden nicht ausreichend

beachtet und befolgt, wie viele Datenskandale aus den letzten Jahren zeigen – wie die NSA-Affäre, die Yahoo-Datenpannen und der Facebook-Cambridge-Analytica-Skandal. Die Organisationen, die private Daten speichern und verarbeiten, z. B. Regierungsbehörden und globale Unternehmen, profitieren davon, so viel wie möglich zu sammeln und zu analysieren, während der Schutz und die Achtung der Privatsphäre eher das Gegenteil erfordert – nur so viel zu erheben und zu speichern, wie implizit erforderlich ist.

Die Datenschutznormen unterscheiden zwischen personenbezogenen Daten und anderen Daten. Personenbezogene Daten sind alle Informationen, die sich auf eine Person beziehen, d. h. alle Informationen, die zur Identifizierung dieser Person verwendet werden können (z. B. Name, Alter, Adresse) sowie alle Informationen, die mit einer Person verknüpft oder verknüpfbar sind (z. B. medizinische Diagnose oder Vorstrafen). Die Datenschutzgesetze legen spezielle Vorschriften für den Schutz der Erhebung, Speicherung und Verarbeitung von personenbezogenen Daten fest, um die Privatsphäre aller Betroffenen zu schützen.

---

### 10.1.2 Wissensschutz

---

Aus Sicht der IT-Sicherheit ist der Schutz des Wissens vergleichbar mit dem Schutz der Privatsphäre. Das Motiv ist jedoch, Wissen vor Wettbewerbern geheim zu halten, um einen Wettbewerbsvorteil zu erlangen und zu erhalten. Der Wissensschutz umfasst den Schutz der Rohdaten selbst, aber auch das Wissen über die Algorithmen und Methoden, die an der Datenanalyse beteiligt sind, sowie das Ergebnis der Analyse.

Der Aspekt des Wissensschutzes unterstreicht das Interesse der Eigentümer von Big-Data-Systemen, ihre Investitionen auf der einen Seite zu schützen. Andererseits haben alle Analyseergebnisse das Potenzial, die Privatsphäre aller zu verletzen, deren Daten zur Berechnung der Ergebnisse verwendet wurden.

Die Interessen derjenigen, die Datenanalysen durchführen, und die der Betroffenen müssen sorgfältig gegeneinander abgewogen werden. Der Bundesgerichtshof hat über die Bewertungsmethoden von Auskunftfeien entschieden, dass es diesen Unternehmen ausreicht, Auskunft darüber zu erteilen, welche Daten gespeichert und zur Berechnung des Kredit-Scores verwendet werden, nicht aber darüber, wie die Analyse im Detail durchgeführt wird [15]. In diesem Fall wird der Schutz des Know-how eines Unternehmens höher gewichtet als das Recht des Einzelnen, zu wissen, welche Daten gespeichert und wie sie verarbeitet werden (in den Gesetzen und der Literatur zum Datenschutz als Transparenzprinzip bezeichnet). Damit ist es möglich, personenbezogene Daten mit einem geheimen Satz von Algorithmen zu verarbeiten, ohne Informationen darüber offenlegen zu müssen, wie die Algorithmen im Detail zu ihren Schlussfolgerungen gekommen sind.

## 10.2 Anforderungen für einen sicheren und datenschutzgerechten Umgang mit Daten

Big Data wird immer innerhalb eines IT-Systems verarbeitet. Daher muss die Sicherheit für Big Data auch die IT-Sicherheit umfassen und während des gesamten Lebenszyklus von Daten in Big-Data-Systemen gewährleistet sein. Ohne eine sichere Umgebung ist eine sichere Erfassung und Verarbeitung

von Daten nicht möglich. Wir geben daher einen kurzen Überblick über die notwendigen Aspekte der IT-Sicherheit, die bei der Nutzung von Big-Data-Systemen nicht zu vernachlässigen sind. Ergänzt wird die Darstellung um besondere Aspekte, die sich aus dem Datenschutz ergeben.

---

### 10.2.1 Übertragung (»Data in Transit«)

---

In Big-Data-Systemen werden Daten kontinuierlich zwischen verschiedenen Geräten, Apps und Servern vom Nutzer zum Betreiber und zurück übertragen. Der Wert der in Big-Data-Systemen gesammelten Daten ergibt sich auch aus der Tatsache, dass die Daten in Echtzeit erfasst und ausgewertet werden. Das Abhören der Kommunikation zwischen Endnutzersystemen und den zentralen Servern bedeutet nicht nur eine Verletzung der Privatsphäre des Nutzers, sondern kann

auch die Vorteile des Dienstbringers einschränken. Dasselbe gilt für die Übertragung zwischen verschiedenen Servern eines Dienstes oder zwischen verschiedenen Diensten. Der Datentransfer zwischen Benutzer-Clients und dem Analysesystem bzw. zwischen verschiedenen Systemen muss vertraulich sein und die Daten müssen während der Übertragung unverändert bleiben.

---

### 10.2.2 Speicherung (»Data at Rest«)

---

Der Umgang mit großen Datenmengen umfasst die Speicherung dieser Daten und darüber hinaus die Speicherung von aggregierten Daten (z. B. Aggregation von Daten zur Modellierung des Verhaltens einer Benutzergruppe) und von Analyseergebnissen. Das Ziel von Organisationen, die mit großen Datenmengen umgehen, ist nur in einer sicheren Umgebung erreichbar. Es ist von grundlegendem Interesse, dass ein unbefugter Zugriff von außen nicht möglich ist und dass der Betreiber in der Lage ist, eine Art Audit und Bewertung durchzuführen, um zu nachzuweisen, dass das System ausreichend sicher ist.

Auf der anderen Seite sollte den Personen, deren Daten aggregiert und gespeichert wurden, die Möglichkeit gegeben werden, zu erfahren, welche Daten über sie tatsächlich gespeichert sind, wie es von den Datenschutzregeln gefordert wird. In der Praxis ist die Erfüllung dieser widersprüchlichen Anforderungen von Zugriffsschutz und Auskunftsmöglichkeit eine wesentliche Herausforderung für Big-Data-Systeme.

---

### 10.2.3 Verarbeitung (»Data in Use«)

---

Die Verarbeitung und Analyse von Daten bleibt in großen Datensystemen weitgehend ein undurchsichtiger Prozess. Einerseits sind Datenverarbeiter nicht in der Lage, die Ergebnisse des Datenanalyseprozesses vorherzusagen, und sie müssen die Analyse möglicherweise im Laufe der Zeit anpassen. Die Analyse selbst ist komplex und kann von Nichtfachleuten kaum vollständig verstanden werden. Gleichzeitig ist das Know-how über die Datenanalyse ein zentraler Bestandteil des Geschäftsmodells von Unternehmen bzw. der Sicherheitsstrategie von Behörden. Die derzeitige Praxis ist, dass Informationen über die Art und Weise, wie die Informationen verarbeitet werden, nicht an die Betroffenen weitergegeben werden, da Organisationen versuchen, ihre eigenen Interessen zu schützen. Die Systeme, die Daten verarbeiten und analysieren, sind daher »Black Boxes«, bei denen der wesentliche Schutz vor allem

durch seine Undurchsichtigkeit gewährleistet ist. Hier gilt es im Einzelfall zu bewerten, welches Maß an Intransparenz mit dem Datenschutz vereinbar ist und was den Betroffenen oder den Aufsichtsbehörden gegenüber offengelegt werden muss.

Während der Verarbeitung ist es schwieriger, die Daten geschützt zu halten, als während der bloßen Speicherung. Besondere Vorsicht ist bei der Auftragsdatenverarbeitung geboten, da dort die Daten sich zur Verarbeitung in fremden Händen befinden. Dies ist in der Praxis häufig der Fall, da der Verarbeiter die nötigen Kapazitäten oder das Know-how und die zur Verarbeitung erforderliche Software hat. Allgemein besteht das Risiko, dass Unbefugte auf das Verarbeitungssystem zugreifen, die Vorgänge beobachten und Daten während der Verarbeitung abgreifen oder verändern.

## 10.3 Ausblick auf die Beiträge dieses Teils der Studie

Big Data erfordert geeignete Schutzmechanismen, um die Daten vor unberechtigtem Zugriff zu schützen; insbesondere dann, wenn es sich um personenbezogene oder -beziehbare Daten handelt. Eine systematische Leitlinie geben die Prinzipien, Konzepte und Strategien zu Privacy by Design aus Kapitel 11. Diese stellen das Bindeglied zwischen den rechtlichen Anforderungen aus Teil II und konkreten technischen Maßnahmen dar. Im Anschluss an Privacy by Design beschreiben wir bestimmte Schutzmaßnahmen, nämlich solche, die aus einer datenzentrischen Sicht auf Big Data besonders wichtig zum Schutz von Privatsphäre sind. Dies sind Schutzmaßnahmen, die direkt auf den Daten arbeiten und dadurch Vertraulichkeit und Integrität sicherstellen. Konkret betrachten wir Methoden zur Verschlüsselung von Daten sowie Methoden zur Anonymisierung von Daten. Dahingegen werden Systemaspekte wie Zugriffskontrolle und Netzwerksicherheit in der Schwesterstudie „Designunterstützung Big-Data-Systeme“ behandelt.

Prinzipiell ist der Einsatz von Verschlüsselungsverfahren sinnvoll, die den Schutz von Data at Rest, Data in Transit und Data in Use gewährleisten können. Da jedoch bei der Auswahl, der Implementierung oder der Anwendung von Verschlüsselung stets eine genaue Abstimmung auf das zu schützende System erforderlich ist, ist zunächst eine Betrachtung der existierenden Verfahren, eine Bewertung ihrer Eignung für bestimmte Anwendungsfälle sowie ein Vergleich untereinander notwendig. In Kapitel 12 wird zunächst ein Überblick über den Stand der Technik eingesetzter Verschlüsselungsverfahren im Kon-

text von Big Data gegeben. Dazu wird separat auf Verfahren für Data at Rest und Data in Transit eingegangen. Neben der Funktionsbeschreibung der verschiedenen Verfahren werden diese hinsichtlich ihrer Eignung und praktischen Anwendbarkeit für den jeweiligen Einsatzzweck verglichen und bewertet. Verschlüsselungsverfahren für Data in Use werden aufgrund der vielen Besonderheiten dieser Verfahren in einem separaten Kapitel erläutert, nämlich in Kapitel 13.

Im Anschluss daran stellen wir verschiedene Konzepte zur Anonymisierung vor. Technische Anonymisierungsverfahren müssen nach der Art der vorliegenden Daten gewählt werden. Wegen der langen Tradition der Verarbeitung und Anonymisierung von strukturierten Daten untersuchen wir zunächst dieses Gebiet (s. Kapitel 14). Da viele Informationen jedoch nicht strukturiert, sondern als Fließtexte vorliegen und zunehmend auch bei automatischen Analysen einbezogen werden, betrachten wir auch die Anonymisierung von Textdaten (s. Kapitel 15). Aufgrund der zunehmenden Relevanz von maschinellem Lernen gehen wir auch auf die Besonderheiten der Anonymisierung in diesem Kontext ein (s. Kapitel 16). Nicht genauer betrachtet wird in dieser Studie die Anonymisierung von Multimediadaten (Bild, Audio, Video), da dies nicht zu den gesteckten Zielen dieser Studie gehört und aufgrund der Weite dieses Feldes eine eigenständige Studie erfordern würde. Verschiedene Herausforderungen und offene Forschungsfragen bei der Anonymisierung von Daten werden in Kapitel 17 dargestellt.

## 11 Privacy by Design und Big Data

In diesem Kapitel werden Konzepte und Ansätze diskutiert, die darauf abzielen, den Datenschutz für Big-Data-Lösungen in einem frühen Stadium zu gewährleisten. Der Begriff Privacy by Design (PbD) wird häufig für IT-Systeme verwendet, die Datenschutzaspekte auf verschiedenen Ebenen behandeln. Dazu gehören sowohl technische Maßnahmen als auch organisatorische Ansätze.

Die große Herausforderung besteht darin, ein System zu entwickeln, das datenschutzfreundlich und dennoch technisch machbar ist. Es ist in der Regel eine Aufgabe, die Kompromisse erfordert: Aus Datenschutzsicht stellt man schnell absolute Anforderungen auf, die technisch kaum zu realisieren sind. Aus Anwendersicht wird die zusätzliche Arbeitsbelastung durch Datenschutzanforderungen oft als Showstopper angesehen. Am Ende muss ein System entstehen, das zumindest den Datenschutzbestimmungen entspricht, aber dennoch attraktiv genug ist, um die Entwicklung durch ein Unternehmen wert zu sein.

Cavoukian hat sieben Grundprinzipien für Privacy by Design eingeführt [16], die in Abbildung 11.1 dargestellt werden. Wir listen sie nachfolgend zusammen mit einer kurzen Erläuterung auf, um die Grundidee von PbD zu veranschaulichen.

**Proaktiv, nicht reaktiv; präventiv, nicht behebend:** PbD versucht, Datenschutzrisiken zu identifizieren, bevor es zu einer Verletzung der Privatsphäre kommt, und schlägt Methoden vor, um deren Auftreten zu verhindern.

**Datenschutzfreundliche Voreinstellungen:** Die Anwendung von PbD bedeutet, dass, wann immer der Benutzer eine Wahl treffen soll, die Standardeinstellung eine die Privatsphäre schützende ist.

**Datenschutz in das Design eingebettet:** Wie der Name schon sagt, verlangt PdD, dass privatheitserhaltende Konzepte und Technologien integraler Bestandteil jedes Systemdesigns sind.

**Volle Funktionalität – Positive Summe, keine Nullsumme:** PbD ist bestrebt, Lösungen anzubieten, die gleichzeitig volle Funktionalität und volle Privatsphäre ermöglichen. Die Privatsphäre sollte nicht als Showstopper oder Overhead angesehen



Abbildung 11.1: Die sieben Grundprinzipien von Privacy by Design.

werden, der nur aufgrund von ComplianceAnforderungen akzeptiert wird.

**EndezuEndeSicherheit – Absicherung des vollständigen Lebenszyklus:** Die EndezuEndeSicherheit gewährleistet, dass es keine Lücken im Schutz privater Daten in Systemen gibt, die nach den Prinzipien von PbD entwickelt wurden. **Sichtbarkeit und Transparenz – Halte es offen:** Ein PbD-System muss Methoden zur Verfügung stellen, um seine Maßnahmen zum Schutz der Privatsphäre zu überprüfen und zu auditieren.

**Respekt vor der Privatsphäre der Benutzer – Fokus auf die Benutzer:** PbD ist ein Konzept zum Schutz der Privatsphäre von Nutzern und Bürgern. Dies sollte immer die wichtigste Perspektive beim Entwurf eines Systems sein.

Diese Prinzipien sind generisch für jedes IT-System. Ihre Einhaltung hilft bei der Entwicklung von Systemen, die naturgemäß die Privatsphäre schützen. In den nächsten Abschnitten werden wir nun einen näheren Blick auf Prinzipien und Konzepte, die speziell für Big Data entwickelt wurden.

## 11.1 Konzepte

Cavoukian und Jonas fassen einige Kernpunkte zusammen, die notwendig sind, um in Big Data privatsphärenhaltende Verarbeitung von personenbezogenen Daten zu ermöglichen [17]. Sie teilen einen hohen Abstraktionsgrad und sind daher eher Gestaltungsprinzipien als Vorschläge für konkrete technische Lösungen.

### **Volle Zuschreibung:**

Das bedeutet, dass jeder Datensatz der gesammelten Daten bis zu seiner Quelle zurückverfolgt werden kann. Während dies für Einzelsätze relativ einfach sein kann, wäre die Aggregation von Daten problematisch, da Daten und Quellen zusammgeführt und gemischt werden. Die volle Zuschreibung erfordert daher eine Art zerstörungsfreie Datenverarbeitung, bei der Operationen, die neue Daten erstellen, den Überblick über die einzelnen beitragenden Datensätze nicht verlieren. Nur dann ist es möglich, bei personenbezogenen Daten die vom Datenschutzgesetz geforderte Löschung oder Korrektur einzelner Datensätze zu ermöglichen.

### **Datenanbindung:**

Wenn Quelldaten geändert werden, sollen die Aktualisierungen auch bei allen Empfängern einer Kopie der Daten aktualisiert werden, so dass deren Analysen nicht auf veralteten Daten stattfinden. Dies ist besonders wichtig bei Daten, die möglicherweise zu Entscheidungen mit negativen Konsequenzen für die Betroffenen führen, z. B. bei Daten für Scoring-Dienste oder für polizeiliche und geheimdienstliche Aktivitäten.

### **Analyse von anonymisierten Daten:**

Wenn Daten geteilt und analysiert werden, kann die Anonymisierung eine sehr hilfreiche Funktion sein. Sie reduziert die mit dem Umgang mit personenbezogenen Daten verbundenen Sicherheitsrisiken weitgehend. Gleichzeitig ermöglichen Big-Data-Analysen in der Regel wertvolle Erkenntnisse auf Basis von Datensätzen, ohne dass personenbezogene Daten benötigt werden.

Die Autoren erwähnen in diesem Zusammenhang »kryptografisch veränderte« Daten, was eher nach Pseudonymisierung als nach Anonymisierung klingt. Dennoch ist der Effekt für den Empfänger der Daten, der sie analysiert, ähnlich.

### **Manipulationssichere Prüfprotokolle:**

Hierbei wird eine Protokolldatei erstellt, in der jeder Zugriff auf

die Datensätze gespeichert wird. Sie ermöglicht zu überprüfen, wer welche Abfrage auf dem Datenbestand ausgeführt hat. Dadurch wird Transparenz über die Datenzugriffe geschaffen und das Vertrauen der Öffentlichkeit in den Big-Data-Analysten erhöht.

Dennoch bleibt eine wichtige Frage: Wer hat Zugriff auf dieses Prüfprotokoll? Solche Protokolle bedeuten auch, dass die Leistung und Arbeitsbelastung von Datenwissenschaftlern überwacht werden kann, was in einigen Umgebungen ein Problem für den Datenschutz und den Arbeitnehmerschutz sein kann.

### **Methoden, die Falsch-Negative bevorzugen:**

In der Regel liefert die Analytik Entscheidungen auf der Grundlage eines Schwellenwerts. Dieser Schwellenwert kann so eingestellt werden, dass entweder Falsch-Negative oder Falsch-Positive wahrscheinlicher sind. Die Autoren schlagen vor, eher Schwellenwerte auszuwählen, die zu mehr Falsch-Negativen führen, insbesondere bei wichtigen Entscheidungen mit großen Auswirkungen auf das Leben des Einzelnen. Durch die Akzeptanz einer höheren Falsch-negativ-Rate reduziert man die Wahrscheinlichkeit, dass z. B. eine falsche Anschuldigung oder Ablehnung auftritt.

### **Selbstkorrigierende Fehlalarme:**

Dieses Konzept fasst einige der früheren Elemente zusammen. Wenn neue Daten abgerufen werden, sollten sie sofort verwendet werden, um bestehende Annahmen zu überprüfen. Wenn ein positives Ergebnis nun zu einem negativen wird, wird es hiermit als ehemaliges Falsch-Positives identifiziert und korrigiert.

Man kann diesen Punkt als eine Kombination aus den strukturellen Anforderungen von Data Tethering und dem Ansatz der Bevorzugung von Falsch-Negativen sehen. Anschließend sollten zunächst positive Ergebnisse oder Vorwürfe überprüft und im Falle einer notwendigen Korrektur so schnell wie möglich kommuniziert werden. Negative, die potenziell positiv werden können, sollten anschließend neu bewertet werden.

### **Informationstransfer-Buchhaltung:**

Es soll ein Protokoll über alle Datenübertragungen erstellt werden. Dies ermöglicht es den Eigentümern des Systems, den Zugriff auf und die Nutzung ihrer Daten zu überprüfen. Diese Anforderung ähnelt dem Prüfprotokoll von oben, bezieht sich aber auf den Datentransfer und nicht auf den Datenzugriff und die Analyse.

## 11.2 Von der Leitlinie zur technischen Umsetzung

Die Europäische Agentur für Netz- und Informationssicherheit<sup>1</sup> (European Union Agency for Network and Information Security, ENISA) hat Designstrategien aus rechtlicher Sicht entwickelt [19]. Sie sind in vier daten- und vier prozessorientierte Strategien unterteilt. Abbildung 11.2 gibt einen Überblick über die Strategien. Neben den Strategien werden auch »Muster« zur Verfügung gestellt, soweit verfügbar, um bei der Umsetzung dieser Strategien zu helfen. Aber diese »Muster« sind oft nur ein Vorschlag für die verfügbare Technologie, ohne echte Ratschläge für deren Umsetzung.

### Minimieren:

Verwenden und speichern Sie so wenig personenbezogene Daten wie möglich. Wenn man dem folgt, kann man sagen, dass jeder Datenverlust oder -missbrauch zumindest die geringstmöglichen Auswirkungen hat. Dies setzt voraus, dass man eine genaue Vorstellung davon hat, welche Art von Analyse mit den Daten durchgeführt werden soll. Es steht in direktem Gegensatz zu der in der Big-Data-Community weitverbreiteten Strategie, zunächst so viele Daten wie möglich zu sammeln und dann herauszufinden, was sich aus den Daten ableiten lässt.

### Verbergen:

Alle personenbezogenen Daten und Mitteilungen sind so zu verschleiern, dass ein menschlicher Beobachter sie nicht verstehen kann. Dies macht den potenziellen Missbrauch durch Angreifer schwieriger und den Verlust personenbezogener Daten durch Vernachlässigung weniger wahrscheinlich. Dies wird in der Regel durch Verschlüsselung, Anonymisierung oder Pseudonymisierung der Daten und durch Verschleierung der Netzwerkkommunikationskanäle erreicht. Während die Verschlüsselung während des Transports oder der Lagerung recht einfach ist, ist die Verschlüsselung während der eigentlichen Verarbeitung der Daten immer noch eine Herausforderung. Für Big Data kann dies auch eine hohe Anforderung sein, da viele Anwendungen aktiv darauf abzielen, mehr über soziale Netzwerke zu erfahren, indem sie die Kommunikation zwischen den Personen analysieren.

### Trennen:

Speichern und verarbeiten Sie personenbezogene Daten einer Person auf verteilte Weise, um die Erstellung vollständiger

Profile dieser Person zu verhindern. Wenn mehrere Datenquellen über eine Person existieren, sollten diese nicht aggregiert werden.

Dies ist natürlich das Gegenteil aller Scoringverfahren, die darauf abzielen, so viele Datenquellen wie möglich über eine Person zu sammeln und zusammenzuführen, um ein detailliertes Profil zu erstellen. Aber auch unabhängig von Anwendungen für das Scoring wird die Aufbewahrung von Daten über eine Person, die über das Netzwerk verteilt sind, die Kosten für die Verarbeitung der Daten erhöhen, sodass es eine Strategie ist, die von den meisten Big-Data-Ansätzen kaum emuliert werden kann.

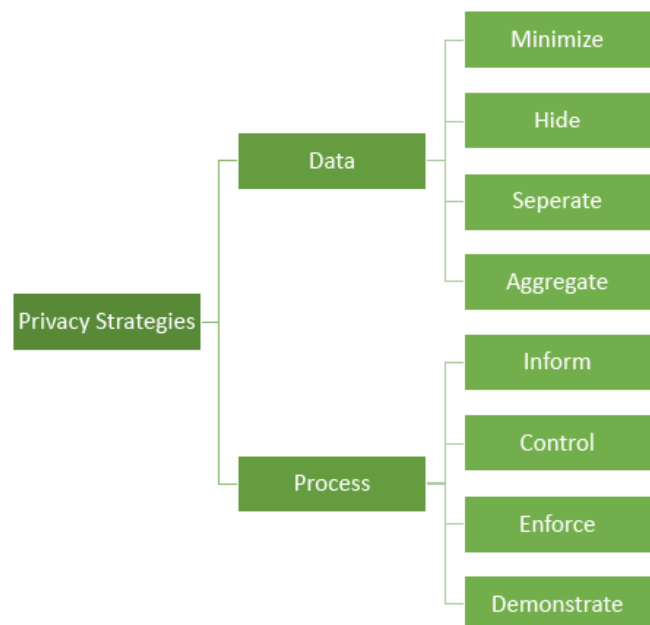


Abbildung 11.2: Die acht Designstrategien.

### Zusammenführen:

Fassen Sie personenbezogene Daten so früh wie möglich in Gruppen zusammen, um die personenbezogenen Daten der Einzelpersonen zu verbergen. Dies kann durch Datenschutztechniken wie die k-Anonymisierung erreicht werden. Obwohl diese Strategie leichter anzuwenden zu sein scheint als andere, kann der erfolgreiche automatisierte Aufbau von Gruppen, die gleichzeitig Privatsphäre und Aussagekraft bieten, eine Herausforderung sein.

<sup>1</sup> [www.enisa.europa.eu](http://www.enisa.europa.eu).



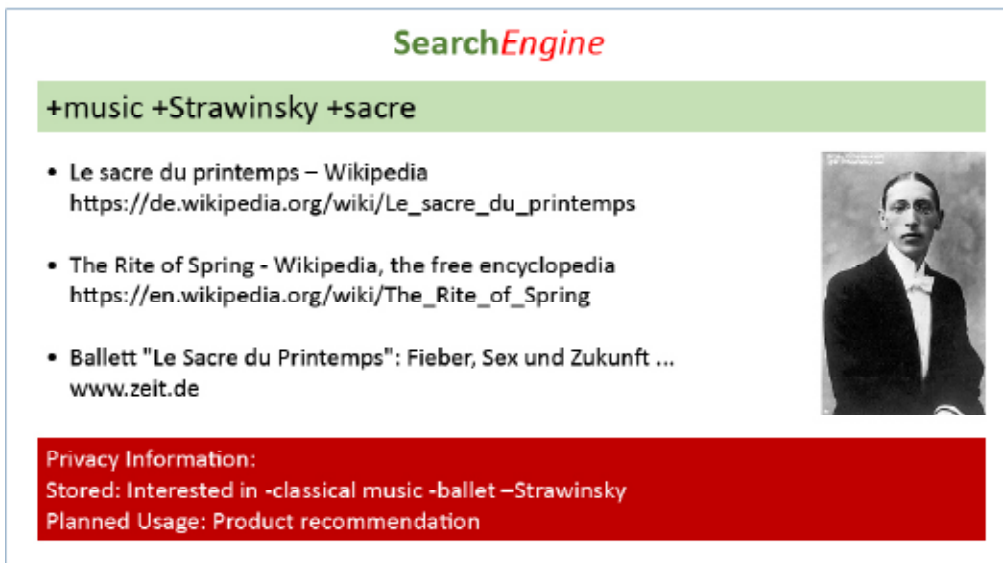


Abbildung 11.3: Beispielmodell zur »Informieren«-Strategie: Die Verwendung einer Suchmaschine erzeugt Daten von Interesse, z. B. für Produktwerbung. Hier wird der Nutzer darüber informiert, dass sein Interesse an klassischer Musik für spätere Anzeigen gespeichert wird.

### Informieren:

Wann immer Daten erhoben werden, sollte man darüber informiert werden, was gespeichert ist, aus welchem Grund die Daten erhoben werden und wie sie geschützt sind. Dies erhöht die Transparenz der Datenerhebung und ermöglicht es den Bürgern, die Zusammenführung der über sie gespeicherten Daten zu verstehen. Die Information über die Weitergabe von Daten an Dritte ist ebenfalls Bestandteil dieser Strategie. Dazu gehört auch, über erfolgreiche Angriffe zu informieren, die zum Verlust privater Daten führen. Um dem Anwender die notwendigen Informationen tatsächlich zur Verfügung zu stellen, sind geeignete Schnittstellen notwendig. Abbildung 11.3 zeigt ein vereinfachtes Beispiel, bei dem gespeicherte Daten und geplante Nutzung am unteren Rand der Schnittstelle erwähnt werden.

### Kontrollieren:

Die Bürger sollten die Kontrolle über die über sie gespeicherten Daten haben. Dies steht in engem Zusammenhang mit der »Informieren«-Strategie, da »Kontrolle« ohne »Informieren« blind wäre und »Informieren« ohne »Kontrolle« frustrierend wäre. Der Grad der Kontrolle ist hier nicht definiert, aber es wird ein Beispiel mit der Möglichkeit gegeben, die Datenschutzeinstellungen in einem sozialen Netzwerk zu ändern. Wie bei der »Informieren«-Strategie ist das Design der Benutzeroberfläche wichtig, um diese Strategie tatsächlich umzusetzen.

Im Rahmen eines Big-Data-Systems würden »Kontrollieren« und »Informieren« eine Schnittstelle und ein Identitätsmanagement für Benutzerkonten erfordern, zum Beispiel als Erweiterung der üblichen Konten. »Informieren« würde Push-Mechanismen benötigen, z. B. per E-Mail als tägliche Zusammenfassung oder direkt bei der Nutzung des Systems. Dies würde zumindest einen erheblichen Aufwand bedeuten, wenn viele Benutzerdaten regelmäßig erhoben werden. »Informieren« könnte möglicherweise zu einem signifikanten Anstieg der Benutzeranfragen führen, die nach weiteren Details fragen oder Erklärungen anfordern, wenn sie diese bereitgestellten Informationen nicht vollständig verstehen.

### Durchsetzen:

Einrichtung und Durchsetzung einer Datenschutzerklärung in Übereinstimmung mit den geltenden Gesetzen. Diese Richtlinie ist eine Kombination aus Regeln und technischen Maßnahmen zu ihrer Durchsetzung. Das notwendige Minimum an Regeln kann aus Gesetzen abgeleitet werden. Technische Maßnahmen können Standardsicherheitsmaßnahmen wie Zugangskontrolle oder Transportverschlüsselung beinhalten. Es wird auch ein System zur Verwaltung von Datenschutzrechten vorgeschlagen, ein System ähnlich einem Digital-Rights-Management-System (DRM-System), das jedoch die Wahrung des Datenschutzes anstelle des Urheberrechts durchsetzt.

**Demonstrieren:**

Man sollte nachweisen können, dass sein System den Datenschutzvorschriften entspricht. Dies erfordert den Nachweis, dass die notwendigen Maßnahmen ergriffen wurden, um die eigenen Datenschutzrichtlinien gegenüber Dritten durchzusetzen. Die Führung von Protokolldateien und die Bereitstellung von Auditierungsmethoden sind dabei von entscheidender Bedeutung. Es bedeutet auch, dass man im Falle einer Ver-

letzung der Privatsphäre schnell die Stelle der Verletzung und die Angriffsstrategie identifizieren kann.

Für komplexe verteilte Systeme, die oft für Big Data erforderlich sind, kann dies eine Herausforderung sein, wenn keine streng durchgesetzte Sicherheitsrichtlinie vorhanden ist. Da IT-Systeme oft ständig um Soft- und Hardware erweitert werden, muss diese Richtlinie regelmäßig aktualisiert werden.

## 11.3 Diskussion

Die oben genannten Prinzipien und Konzepte helfen sicherlich, sich ein Big-Data-System vorzustellen, das darauf abzielt, die Analyse personenbezogener Daten und den Datenschutz gleichzeitig zu ermöglichen. Daher sind sie sehr hilfreich in der Designphase einer Big-Data-Lösung.

Aus der Sicht eines Praktikers sind sie jedoch noch recht theoretisch: Zum Beispiel ist es ein erster Schritt zu wissen, dass ein vollständiges Protokoll aller Datenzugriffe und -bewegungen hilfreich wäre, damit ein System einfach zu auditieren ist, aber es ist noch ein langer Weg, um einen solchen Mechanismus tatsächlich zu erreichen. Und seine Auswirkungen auf das Gesamtsystem bleiben unklar. Zumindest besteht die Gefahr, dass die detaillierte Überwachung komplexer wird als das Kernsystem.

Dies ist eine häufige Kritik bei den Datenschutzerfordernissen: Sie sind leicht zu nennen, aber schwer zu erfüllen. Aber

diese Denkweise kann hauptsächlich von den Erfahrungen beeinflusst werden, die man gemacht hat, indem man bereits bestehende Lösungen um privatsphärenhaltende Mechanismen erweitert hat. Dies ist in der Regel eine sehr teure und komplexe Aufgabe. Aus diesem Grund kann der konstruktive Aspekt in diesem Kapitel nicht ausreichend betont werden. Die Neukonzeption von Systemen mit den oben genannten Konzepten wird zu Systemen führen, bei denen die Integration von Mechanismen zur Wahrung der Privatsphäre wesentlich einfacher wird. Am Ende wird der Datenschutz durch Design für große Datenmengen ein evolutionärer Prozess sein, der das Systemdesign in den folgenden Jahren beeinflusst.

Die nächsten Kapitel werden einen Überblick über konkretere Mechanismen geben, um die Ziele, die die Prinzipien und Konzepte aus diesem Kapitel anregen, tatsächlich zu erreichen.

# 12 Verschlüsselungsmechanismen für Daten in Big-Data-Systemen

Wesentliche Methoden zur sicheren Datenspeicherung und -übertragung sind Verschlüsselungsmechanismen. Während Verschlüsselung bei Webbrowsern, Instant Messengern, E-Mail und vielen weiteren alltäglichen Anwendungen bereits selbstverständlich oder gängig ist, wird Verschlüsselung insbesondere bei Big Data noch sehr argwöhnisch betrachtet. Das liegt zum einen an fehlendem Know-how über Verschlüsselungstechnologie, zum anderen daran, dass die zusätzlichen Operationen (Ver- und Entschlüsseln, Schlüsselaustausch und -verwaltung) die Performanz mindern und zusätzliche Konfigurations- und Verwaltungsarbeit verursachen.

Auch im Big-Data-Umfeld bildet Kryptografie den Kern von

Schutzmechanismen für sämtliche geschäftlich schützenswerte Daten – insbesondere für personenbezogene und -beziehbare Daten –, die im Kontext von Big Data gespeichert, ausgetauscht und verarbeitet werden. Nach der Erklärung kryptografischer Grundbegriffe werden die drei Zustände »Data at Rest«, »Data in Transit« und »Data in Use« mit ihren kryptografischen Anforderungen beschrieben. In den nachfolgenden Abschnitten werden für »Data at Rest« und »Data in Transit« geeignete und gängige kryptografische Konzepte und Verfahren beschrieben. Kryptografische Konzepte für »Data in Use« werden aufgrund ihrer gänzlich anderen Natur und aufgrund der Komplexität des Themas in Kapitel 13 beschrieben.

## 12.1 Grundbegriffe

Zunächst folgt eine kurze Definition der wichtigsten Begriffe aus dem Bereich Kryptografie, die als Grundlage für die nachfolgenden Abschnitte dient.

### **Schlüssel (Key):**

Ein Schlüssel ist eine Information, die zur Verschlüsselung bzw. Entschlüsselung verwendet wird. Hierbei ist zu beachten, dass ein Schlüssel ausdrücklich nicht äquivalent zu einem Passwort ist: Letzteres wird selten direkt zur Ver-/Entschlüsselung eingesetzt, sondern meist wird der eigentliche Schlüssel vom Passwort abgeleitet oder der Schlüsselzugriff wird über das Passwort kontrolliert.

### **Verschlüsselungsverfahren (Cipher):**

Ein Verschlüsselungsverfahren bezeichnet den eigentlichen Algorithmus, der zur Verschlüsselung bzw. Entschlüsselung verwendet wird. Jedes Verschlüsselungsverfahren benötigt mindestens einen Klartext (Plaintext) und einen Schlüssel, um diesen Klartext zu einem Geheimtext, genannt Chiffre (Ciphertext), zu verschlüsseln, und es benötigt für die Entschlüsselung mindestens ein Chiffre und einen Schlüssel, um wieder den zugehörigen Klartext aus dem Chiffre zu erzeugen.

### **Symmetrische Verschlüsselungsverfahren:**

Verwendet ein Verschlüsselungsverfahren denselben Schlüssel zur Verschlüsselung wie zur Entschlüsselung, spricht man von einem symmetrischen Verfahren.

### **Asymmetrische Verschlüsselungsverfahren:**

Werden bei einem Verschlüsselungsverfahren für die Verschlüsselung und für die Entschlüsselung jeweils verschiedene Schlüssel verwendet, spricht man von einem asymmetrischen Verschlüsselungsverfahren. Die beiden verschiedenen Schlüssel bilden dabei ein Schlüsselpaar (Key pair), das aus einem öffentlichen Schlüssel (Public key) und einem privaten Schlüssel (Private key) besteht, wobei der private Schlüssel nicht mit realisiertem Aufwand mithilfe des öffentlichen Schlüssels oder auch mithilfe von Signaturen, die mit dem privaten Schlüssel erstellt wurden, bestimmt werden kann.

### **Privater Schlüssel (Private key):**

Bei asymmetrischen Verfahren dient der private Schlüssel der Entschlüsselung von Informationen, die mit dem zugehörigen öffentlichen Schlüssel verschlüsselt wurden. Weiterhin wird der private Schlüssel zur Anfertigung von digitalen Signaturen eingesetzt.

**Öffentlicher Schlüssel (Public key):**

Bei asymmetrischen Verfahren dient der öffentliche Schlüssel zur Verschlüsselung von Informationen. Weiterhin wird der öffentliche Schlüssel zur Überprüfung von digitalen Signaturen, die mit dem zugehörigen privaten Schlüssel angefertigt wurden, eingesetzt.

**Digitale Signatur:**

Eine digitale Signatur ist ein digitales Objekt, welches der Verfasser oder Absender eines Dokumentes mittels eines asymmetrischen kryptografischen Verfahrens aus dem Dokument unter Verwendung seines privaten Schlüssels erzeugt und welches das Dokument mit Integritätsschutz sowie einseitigem Schutz vor Abstreitbarkeit versieht. Weiterhin kann mit einer digitalen Signatur die Authentizität eines Dokumentes zweifelsfrei festgestellt werden, wenn der zum privaten Schlüssel des Absenders zugehörige öffentliche Schlüssel eindeutig dem Absender zugeordnet werden kann.

**Digitales Zertifikat:**

Ein digitales Zertifikat ist ein digitales Objekt, das Eigenschaften von anderen Objekten oder Personen mittels asymmetrischer Verschlüsselung bestätigen kann. Beispielsweise lässt sich durch ein Public-Key-Zertifikat die Identität des Besitzers des im Zertifikat enthaltenen öffentlichen Schlüssels bzw. des entsprechenden privaten Schlüssels feststellen. Umgesetzt wird die Bestätigung der in digitalen Zertifikaten enthaltenen Eigenschaften mithilfe einer Zertifizierungsstelle (Certificate Authority, CA): Diese stellt das digitale Zertifikat aus und versichert dabei mit einer digitalen Signatur, die im Zertifikat enthaltenen Eigenschaften überprüft und für korrekt befunden zu haben (insb. Angaben über den Besitzer) oder so festgelegt zu haben (bspw. die Gültigkeitsdauer). Mit der Gültigkeit dieser digitalen Signatur werden dann die Eigenschaften als bestätigt angesehen. Hierfür notwendig ist ein Grundvertrauen in eine Reihe von Zertifizierungsstellen, das in Form von mitgelieferten Public-Key-Zertifikaten dieser Zertifizierungsstellen (sogenannten Root-Zertifikaten) sowohl auf Betriebssystemebene als auch auf Applikationsebene (bspw. im Webbrowser) hergestellt wird.

**Blockchiffre (Block Cipher):**

Eine Blockchiffre ist ein Verschlüsselungsverfahren, das auf einen Klartext fester Länge (Klartextblock) angewandt werden kann und diesen auf einen verschlüsselten Block (Ciphertextblock) abbildet. Dies ist die natürliche Arbeitsweise der meisten Verschlüsselungsalgorithmen.

Im Gegensatz zu Stromchiffren können Blockchiffren nur einen einzigen Klartextblock verschlüsseln. Daher müssen Blockchiffren bei einer Anwendung auf einen Klartext mit einer Länge, die größer als die vorgegebene Blocklänge ist, einen sogenannten Betriebsmodus (Mode of Operation) wählen, der festlegt, wie die Blockchiffre auf mehrere Blöcke angewandt wird (siehe Abschnitt 12.2).

**Stromchiffre (Stream Cipher):**

Eine Stromchiffre ist ein Verschlüsselungsverfahren, das auf einen Klartext beliebiger Länge angewandt werden kann. Hierbei wird zunächst aus einem zufälligen Startwert (Initialisierungsvektor) und dem kryptografischen Schlüssel ein pseudzufälliger Schlüsselstrom erzeugt, der anschließend mit dem zu verschlüsselnden Klartext XOR-verknüpft wird. Für Stromchiffren ist es im Gegensatz zu Blockchiffren nicht notwendig, dass bereits ein Klartext(block) vollständig vorliegt, bevor eine Verschlüsselung durchgeführt werden kann. Für Stromchiffren ist von besonderer Relevanz, dass der Initialisierungsvektor niemals ein zweites Mal eingesetzt wird.

**Kryptografischer Hash:**

Ein kryptografischer Hash ist das Ergebnis einer kryptografischen Hashfunktion, bei der eine beliebige Information als Eingabe dient, welche auf eine Information mit fester Länge (Hash) abgebildet wird. Im Vergleich zu einem gewöhnlichen Hash haben kryptografische Hashes zusätzliche Sicherheitsmerkmale, die es z. B. quasi unmöglich machen, für einen gegebenen kryptografischen Hash einen Eingabewert zu finden – weder den Original-Eingabewert noch einen anderen Eingabewert –, der auf denselben Hash abgebildet wird (Urbild-Resistenz). Ebenso besteht praktisch keine Chance, für einen gegebenen Eingabewert einen anderen, verschiedenen Eingabewert zu finden, der sich auf denselben kryptografischen Hash abbilden lässt (Zweites-Urbild-Resistenz). Zusätzlich ist bei sicheren Hashfunktionen auch die Suche nach zwei unterschiedlichen Eingabewerten, die über die kryptografische Hashfunktion auf denselben Hash abgebildet werden, praktisch unmöglich (Kollisionsresistenz).

## 12.2 Betriebsmodi für Blockchiffren

Blockchiffren wie AES werden mit einem Betriebsmodus derart erweitert, dass sie auch Sequenzen von Blöcken verschlüsseln können und teilweise sogar zu Stromchiffren werden. Nachfolgend wird für eine Reihe von Betriebsmodi auf bestehende Vor- und Nachteile hingewiesen:

### ECB:

»Electronic Codebook« (ECB) beschreibt einen sehr einfachen Betriebsmodus, bei dem ein Klartext beliebiger Länge in Blöcke mit einer fixen Länge aufgeteilt wird, die jeweils einzeln verschlüsselt werden. Die Entschlüsselung wird analog durchgeführt: Der Ciphertext wird erneut in Blöcke aufgeteilt, die jeweils einzeln entschlüsselt werden und nach der Entschlüsselung in ihrer Gesamtheit den Klartext ergeben.

Ein sehr großer Nachteil bei der Verwendung von ECB ist, dass einzelne Blöcke, die identischen Klartext enthalten, nach der Verschlüsselung ebenfalls den einen identischen Ciphertext ergeben. Somit sind Muster im Klartext auch nach der Verschlüsselung identifizierbar. Von der Verwendung von ECB wird somit gänzlich abgeraten.

### CBC:

Der Betriebsmodus »Cipher Block Chaining« (CBC) behebt Nachteile von ECB, indem jeder Klartextblock mit dem vorangehenden Ciphertextblock XOR-verknüpft wird, bevor die eigentliche Verschlüsselung des Blocks stattfindet. Der erste Klartextblock wird entsprechend mit einem zufällig gewählten Initialisierungsvektor (IV) XOR-verknüpft. CBC erlaubt aufgrund der XOR-Verknüpfung zwar keine parallelisierte Verschlüsselung, die Entschlüsselung kann jedoch parallel ausgeführt werden, da bereits alle verschlüsselten Blöcke vorliegen. Weiterhin muss bei CBC die Länge des Klartexts einem Vielfachen der verwendeten Blockgröße entsprechen, weshalb der Klartext ggf. mit weiteren fest definierten Bytes aufgefüllt wird (Padding).

CBC ist anfällig für sogenannte »Padding Oracle Attacks«, mit denen eine Entschlüsselung von Ciphertext möglich ist, der mit einem Blockchiffre unter Verwendung von CBC verschlüsselt wurde. Dies geschieht dadurch, dass Fehlerfälle, die das eingesetzte Padding bei CBC betreffen, bei einer Applikation erzeugt werden, die für die Entschlüsselung von Ciphertext im CBC-Modus verwendet wird. Anschließend wird das Verhalten der Applikation in diesen Fehlerfällen ausgenutzt, um Byte für Byte den Klartext abzuleiten. Der Verschlüsselungsalgorithmus

selbst ist dabei unerheblich, da die Attacke explizit auf Eigenschaften des Betriebsmodus CBC abzielt. Ob eine Applikation, die Blockchiffren mit CBC einsetzt, anfällig für Padding Oracle Attacks ist, hängt von der jeweiligen Implementierung der Applikation ab. Bekannte Beispiele für erfolgreiche Padding Oracle Attacks gegen SSL und TLS sind POODLE [58] und Lucky Thirteen [2].

### CTR:

Der Betriebsmodus »Counter« (CTR) wandelt einen Blockchiffre effektiv in einen Stromchiffre um, indem zunächst für jeden Klartextblock ein Schlüsselstromblock erzeugt wird, der anschließend mit dem jeweiligen Klartextblock XOR-verknüpft wird. Der Schlüsselstromblock wird hierbei durch das Verschlüsseln eines Wertes erzeugt, der meist aus einem ansteigenden Zahlenwert und einer zufälligen Nonce bzw. einem Initialisierungsvektor gebildet wird. Dieser Wert wird namensgebend für diesen Betriebsmodus als »Counter« bezeichnet. Dadurch, dass beim Verschlüsseln statt des vorangehenden Ciphertext-Blocks ein Counter verwendet wird, ist sowohl das Verschlüsseln als auch das Entschlüsseln im Betriebsmodus CTR sehr gut parallelisierbar.

### OCB:

Der »Offset Codebook Mode« (OCB) verwendet ebenfalls wie der CTR-Modus einen Counter, fügt jedoch noch einen Message Authentication Code (MAC) hinzu, der zur Überprüfung der Integrität und Authentizität der verschlüsselten Daten verwendet werden kann. OCB arbeitet weiterhin sehr effizient, da nur eine Operation pro Block benötigt wird, um einen Ciphertextblock mit Authentifizierungsdaten zu erzeugen. OCB unterliegt US-Patenten, sodass die Nutzung innerhalb der USA nur eingeschränkt möglich ist.

### GCM:

Ein weiterer, auf einem Counter basierender Betriebsmodus ist der »Galois/Counter Mode« (GCM). Das Grundprinzip entspricht dem des CTR- und OCB-Modus, lediglich die Authentifizierungsdaten werden im Unterschied zu OCB durch Multiplikationen im Galois-Körper  $GF(2^{128})$  erzeugt. Wie auch OCB und CTR lässt sich GCM sehr gut parallelisieren und kann außerdem bei aktuellen CPUs von einer speziellen Befehlsweiterung (CLMUL) profitieren, die sich u. a. zur effizienten Implementierung von Multiplikationen im Galois-Körper eignen.

GCM wird weitläufig als Standardbetriebsmodus für symmetrische Verschlüsselungsverfahren empfohlen.

**XTS:**

Der Betriebsmodus »XEX Tweakable Block Cipher with Cipher-text Stealing« (XTS) [31] ist speziell für die Verschlüsselung

von Inhalten geeignet, auf die ein wahlfreier Zugriff (engl. Random Access), d. h. ein Zugriff an beliebiger Stelle, möglich sein soll (bspw. Festplatteninhalte). XTS ist für diesen Einsatzzweck sehr gut geeignet und sollte aufgrund fehlender Authentifizierungsmechanismen ausschließlich für diesen Zweck verwendet werden.

## 12.3 Zustände digitaler Daten in Big-Data-Systemen

Für den Schutz digitaler Daten betrachtet man drei verschiedene Zustände, in denen sich Daten befinden können. Für jeden dieser Zustände existieren zum Schutz der jeweiligen

Daten bestimmte Anforderungen, die sich durchaus erheblich voneinander unterscheiden.

---

### 12.3.1 Data at Rest

---

Mit »Data at Rest« wird der Zustand innerhalb eines Big-Data-Systems beschrieben, bei dem Daten persistent auf einem Speichermedium gehalten werden, sich also momentan nicht in der Verarbeitung oder im Transfer befinden. Zum Schutz von Data at Rest eignen sich symmetrische Verschlüsselungsverfahren insbesondere wegen ihrer hohen Performanz gegenüber asymmetrischen Verschlüsselungsverfahren. Im Vergleich zu Data in Transit ist bei der Verschlüsselung von

Data at Rest kein Schlüsselaustauschverfahren notwendig. Der symmetrische Schlüssel wird stattdessen einmalig für den Zugriff auf die entschlüsselten Daten in den Speicher geladen und bleibt dort solange vorhanden, bis kein Zugriff mehr auf die entschlüsselten Daten benötigt wird. Entsprechend ist die Eingabe eines Passworts üblicherweise nur beim Start des gesamten Systems notwendig.

---

### 12.3.2 Data in Transit

---

»Data in Transit« (oft auch als »Data in Motion« bezeichnet) beschreibt den Zustand, bei dem Daten zwischen zwei Kommunikationspartnern transferiert werden, sich also aktuell weder in der Verarbeitung (»Data in Use«) noch im abgespeicherten Zustand auf einem Datenträger (»Data at Rest«) befinden. Für die Verschlüsselung von »Data in Transit« ist neben der Verschlüsselung der eigentlichen Nutzdaten auch die Authentizität der Daten wichtig: Jedem Kommunikationsteilnehmer muss über geeignete Verfahren die Möglichkeit

gegeben werden, die Identität seines Gegenübers zweifelsfrei feststellen zu können, um beispielsweise Angriffe, bei denen sich ein Angreifer als der jeweils andere Kommunikationspartner ausgibt (Man in the Middle (MITM) Attacks), abzuwehren. Entsprechend ist neben einer symmetrischen Verschlüsselung der Daten auch die Verwendung einer Public-Key-Infrastruktur (asymmetrische Verschlüsselung) notwendig, mittels derer ein authentifizierter Schlüsselaustausch durchgeführt werden kann.

---

### 12.3.3 Data in Use

---

»Data in Use« beschreibt den Zustand, in dem sich Daten innerhalb eines Big-Data-Systems befinden, wenn diese augenblicklich verarbeitet werden, indem beispielsweise eine Analyse darauf durchgeführt wird. Die Besonderheit beim Schutz von »Data in Use« durch Verschlüsselung liegt darin,

dass zunächst eine Verarbeitung von verschlüsselten Daten ohne Kenntnis des jeweiligen Schlüssels intuitiv nicht möglich erscheint. Tatsächlich existieren jedoch Ansätze, um bestimmte Verarbeitungsschritte auf speziell verschlüsselten Daten auszuführen, ohne dass diese entschlüsselt werden müssen.

## 12.4 Geeignete Verschlüsselung für »Data at Rest«

Seit mehr als 15 Jahren ist der Advanced Encryption Standard (AES) das Mittel der Wahl für symmetrische Verschlüsselung. Zwar existieren auch zahlreiche weitere geeignete symmetrische Verschlüsselungsverfahren, jedoch ist AES mittlerweile einer der am weitesten verbreiteten Standards und kann u. a. mittels eines eigenen Befehlssatzes für CPUs (AES-NI) effizienter verwendet werden, wodurch noch zusätzliche Geschwindigkeitssteigerungen mit AES ermöglicht werden.

AES kann mit den Schlüssellängen 128 Bit, 192 Bit und 256 Bit eingesetzt werden. Bisher bekannte Key Recovery Attacks auf AES erreichen eine Reduzierung der Schlüssellänge um maximal 2 Bit, sodass AES selbst unter Verwendung eines 128-Bit-Schlüssels aktuell als sehr sicher zu bezeichnen ist.

Obwohl sich AES als Verschlüsselungsverfahren für »Data at Rest« sehr gut eignet, ist es wichtig zu wissen, welche Betriebsmodi für Blockchiffren wie AES eingesetzt werden können, ohne ein Sicherheitsrisiko zu erzeugen, vgl. Abschnitt 12.2.

Möchte man Verschlüsselung einsetzen, ist ein intelligentes Schlüsselmanagement unabdingbar. Hier liegt auch die große Herausforderung für Big-Data-Systeme. Die Installation und Instandhaltung sowie die Sicherung eines solchen Systems ist zusätzliche Arbeit bzw. verursacht zusätzliche Kosten für einen erst einmal nicht erfassbaren Nutzen. Hier kann beispielsweise Apache Ranger<sup>1</sup> die Arbeit erleichtern.

## 12.5 Geeignete Verschlüsselung für »Data in Transit«

Für die Transportverschlüsselung (»Data in Transit«) in Big-Data-Systemen existiert mit der Transport Layer Security (TLS) ein leistungsfähiges Verschlüsselungsprotokoll, das auf Basis moderner Kryptografie Vertraulichkeit und Authentizität für die zu übertragenen Daten liefert. TLS ist der Nachfolger von Secure Socket Layer (SSL), das bis inklusive Version 3.0 als veraltet gilt und auch aufgrund zahlreicher Sicherheitslücken nicht mehr eingesetzt werden sollte.

Um Authentizität zu gewährleisten, wird bei TLS eine zertifikatsbasierte Authentifizierung von Kommunikationsteilnehmern eingesetzt.<sup>2</sup> Meist wird nur die Authentifizierung eines Servers durch den jeweiligen Client durchgeführt, wobei der Client die Gültigkeit des Server-Zertifikats überprüft. Ebenfalls ist durch die Verwendung von Client-Zertifikaten auch die Authentifizierung von Clients gegenüber dem Server und anderen Clients möglich. Bei der Verwendung von TLS sind neben der Wahl der Zertifizierungsstelle (Certificate Authority, CA), die für das jeweilige Big-Data-System gültige Zertifikate ausstellt, eine Reihe von Punkten zu beachten, die maßgeblich zur Sicherheit des Protokolls beitragen. Nachfolgend werden diese Punkte mit den entsprechenden Empfehlungen aufgelistet.

### Protokollversionen:

Die Vorgängerversionen von TLS (SSL 2.0, SSL 3.0) sind aufgrund zahlreicher teilweise sehr kritischer Sicherheitsrisiken als veraltet zu bezeichnen und sollten heute nicht mehr eingesetzt werden [6, 58]. TLS 1.0 enthält ebenfalls Schwachstellen [26], wird aber in der Praxis noch von zeitgemäßen Servern und Clients unterstützt, da nach wie vor verbreitete älteren Clients und schlecht konfigurierte Server keine neueren Versionen unterstützen. Optimalerweise sollte heute ausschließlich TLS 1.1 und 1.2 eingesetzt werden. Letztere Version unterstützt insbesondere moderne Verschlüsselungsverfahren. Aktuell wird an der Entwicklung von TLS 1.3 gearbeitet, mit der weitere neue Verschlüsselungsverfahren in TLS unterstützt und veraltete Verfahren verworfen werden.

### Schlüssellängen:

Für die bei TLS eingesetzten Zertifikate sollten die zugehörigen privaten Schlüssel mindestens eine Schlüssellänge von 2048 Bit bei der Verwendung von RSA-Schlüsseln bzw. 256 Bit bei der Verwendung von ECDSA-Schlüsseln, die auf Kryptografie mit elliptischen Kurven basieren, aufweisen. Kürzere Schlüssellängen bieten nach dem heutigen Stand der Technik keinen

<sup>1</sup> Dies und weitere Möglichkeiten zur Zugriffsverwaltung werden in der Schwesterstudie »Designunterstützung für Big-Data-Systeme« im Kapitel zu Zugriffskontrolle behandelt.

<sup>2</sup> Weitere Informationen hierzu sind in der Schwesterstudie »Designunterstützung für Big-Data-Systeme« unter dem Thema »authentifizierten Schlüsselaustausch« im Kapitel zu Zugriffskontrolle zu finden.

ausreichenden Schutz gegen Angriffe wie die Faktorisierung von RSA-Modulen. Weitere Empfehlungen für Schlüssellängen sind vom Bundesamt für Sicherheit in der Informationstechnik (BSI) verfügbar [14].

### Kryptografische Hashes:

Bei TLS werden u. a. für die Überprüfung von Zertifikaten kryptografische Hashes eingesetzt, die als Fingerabdruck dienen: Stimmt der Fingerabdruck eines Zertifikats mit dem bekannten Fingerabdruck überein, kann mit an Sicherheit grenzender Wahrscheinlichkeit angenommen werden, dass das Zertifikat seit der letzten Überprüfung nicht verändert wurde. Damit bestätigt eine digitale Signatur, die auf dem Hash berechnet wurde, die Echtheit des gesamten Zertifikats, d. h. die Zertifizierungsstelle, welche die Signatur erstellt hat, hat dieses Zertifikat so herausgegeben. Vertraut man der Zertifizierungsstelle, dann kann man somit auch dem Zertifikat vertrauen, selbst wenn man es zuvor nicht gesehen hat. Dies ist das bei TLS übliche Vertrauensprinzip.

Kryptografische Hashfunktionen werden in TLS ebenso in der Funktion HMAC (»Hash-based Message Authentication Code« bzw. »Keyed-Hash Message Authentication Code«) verwendet, welche sowohl zur Berechnung von »Message Authentication Codes« (MACs) für die Prüfung der Authentizität der übertragenen Daten als auch als Baustein einer Pseudozufallsfunktion (Pseudo Random Function, PRF) eingesetzt wird. Die PRF wird für die Vereinbarung eines gemeinsamen Geheimnisses während des Sitzungsaufbaus und für weitere Operationen verwendet.

Es sollte insbesondere von der Verwendung von kryptografischen Hashfunktionen wie MD5 oder SHA-1 abgesehen werden: Diese gelten als veraltet und weisen teilweise signifikante Sicherheitsmängel auf, sodass von einer weiteren Benutzung in TLS abgeraten wird. Moderne Clients (insb. gängige Webbrowser in den Anfang 2017 veröffentlichten Versionen) geben eine Sicherheitswarnung bei SHA-1-Zertifikaten; MD5-Zertifikate werden bereits seit längerer Zeit nicht mehr akzeptiert. Zur Nachrichtenauthentifizierung akzeptieren die Clients noch die Verwendung von SHA-1.

### Verschlüsselungsverfahren:

TLS unterstützt eine Reihe von kryptografischen Verfahren, die für den Schlüsselaustausch und die spätere Verschlüsselung

der Nutzdaten eingesetzt werden können. In der Vergangenheit sind für eine ganze Reihe dieser Verschlüsselungsverfahren teilweise schwerwiegende Sicherheitsmängel aufgedeckt worden, die so gravierend sind, dass vom Einsatz dieser Verfahren gänzlich abgeraten bzw. der Einsatz im Rahmen von Standardisierungsdokumenten wie diversen RFCs verboten wurde. Ein Beispiel für ein Verschlüsselungsverfahren, das heute bei der Verwendung von TLS nicht mehr eingesetzt werden sollte, ist RC4 [3]. Nach RFC 7465 ist der Einsatz von RC4 in TLS mittlerweile verboten [63].

Bei TLS kann über den Konfigurationsparameter »Cipher Suites« festgelegt werden, welche Verschlüsselungsverfahren sowie Verfahren für Schlüsselaustausch oder MACs mit welcher Präferenz eingesetzt werden sollen. Für den typischen Einsatz von TLS auf einem Webserver werden beispielsweise von der Mozilla Foundation unterschiedliche Konfigurationen von Cipher Suites empfohlen, die je nach Anwendungsfall auf möglichst modernen Verschlüsselungsverfahren oder auf hohe Kompatibilität ohne signifikante Defizite bezüglich der Sicherheit ausgerichtet sind.

Folgende Liste von Cipher Suites gemäß der »Modern«-Konfiguration<sup>3</sup> von Mozilla setzt ausschließlich Verfahren aus TLS 1.2 ein:

- ECDHE-ECDSA-AES256-GCM-SHA384
- ECDHE-RSA-AES256-GCM-SHA384
- ECDHE-ECDSA-CHACHA20-POLY1305
- ECDHE-RSA-CHACHA20-POLY1305
- ECDHE-ECDSA-AES128-GCM-SHA256
- ECDHE-RSA-AES128-GCM-SHA256
- ECDHE-ECDSA-AES256-SHA384
- ECDHE-RSA-AES256-SHA384
- ECDHE-ECDSA-AES128-SHA256
- ECDHE-RSA-AES128-SHA256

<sup>3</sup> [https://wiki.mozilla.org/Security/Server\\_Side\\_TLS#Modern\\_compatibility](https://wiki.mozilla.org/Security/Server_Side_TLS#Modern_compatibility).



Kryptografie auf elliptischen Kurven wird hier für den Schlüsselaustausch erzwungen und für Zertifikate bevorzugt. SHA-1 wird als kryptografischer Hash nicht mehr eingesetzt, ebenfalls fehlt die Unterstützung für die mit Sicherheitsproblemen belasteten Verschlüsselungsverfahren RC4 und 3DES [3, 8].

**Zertifizierungsstellen:**

Bei der Wahl der CA sind die existierenden Sicherheitsvorkehrungen für die CA maßgeblich, da diese befugt ist, Zertifikate auszustellen, die von jeder Instanz eines Big-Data-Systems, die

TLS verwendet, als gültig betrachtet werden. Entsprechend stellt eine CA ein äußerst vielversprechendes Ziel für Angreifer dar. CAs sollten regelmäßig Sicherheitsaudits unterzogen werden, bei denen neben den Mechanismen, die für den Schutz des privaten Schlüssels des Root-Zertifikats eingesetzt werden, auch beispielsweise der Umgang mit Sicherheitsvorfällen in der Vergangenheit untersucht wird. Dies ist insbesondere für CAs relevant, die nur innerhalb des Unternehmens verwendet werden.

## 13 Geeignete Verschlüsselung für Data in Use

Die Verschlüsselung von Daten ist heute verbreitet, wenn die Daten transportiert oder gespeichert werden. Demgegenüber ist die Verschlüsselung von Daten, während sie sich in Verwendung befinden, noch eine Herausforderung. Wenn Daten verarbeitet werden, befinden sie sich in der Regel in einem unverschlüsselten Zustand. Ein Angreifer, der an dieser Stelle lauschen kann, hat Zugriff auf die Daten. In einer Big-

Data-Anwendung, die auf Cloud Computing basiert, könnte zumindest der Cloud-Anbieter auf die Daten zugreifen. Homomorphe Verschlüsselung (Abschnitt 13.1) und sichere Mehrparteienberechnung (Abschnitt 13.2) sollen diesem Sicherheitsrisiko entgegenwirken: Im Gegensatz zur herkömmlichen Verschlüsselung können Daten in einem verschlüsselten Zustand verarbeitet werden.

### 13.1 Homomorphe Verschlüsselung

Homomorphe Verschlüsselung hat das Ziel, auf verschlüsselten Daten arbeiten zu können. Homomorph verschlüsselte Daten werden zur Verarbeitung nicht entschlüsselt und auch die Berechnungsergebnisse liegen direkt in verschlüsselter Form vor.

Homomorphe Verschlüsselung bietet neue Sicherheitsgarantien insbesondere für Anwendungen wie Cloud Computing: Die Daten werden vom legitimen und vertrauenswürdigen Inhaber der Daten verschlüsselt, an einen nicht vertrauenswürdigen Dienst gesendet, dort verarbeitet, während sie in einem verschlüsselten Zustand bleiben, und die abgeleiteten Ergebnisse sind ebenfalls verschlüsselt. Erst nachdem die Ergebnisse wieder in eine sichere Umgebung übertragen wurden, werden sie entschlüsselt und verwendet. Dieser Prozess ist in Abbildung 13.1 dargestellt.

Neben Cloud Computing kann die homomorphe Verschlüsselung als Baustein in kryptografischen Protokollen für verschiedene Anwendungsszenarien eingesetzt werden, darunter privatsphärenhaltende Suchmaschinen, Spam-Filterung verschlüsselter E-Mails, Softwareschutz, Wiederverschlüsselung von Daten unter einem neuen Schlüssel ohne Entschlüsselung, Blockchains mit starken Datenschutzgarantien, elektronische Abstimmungen und sichere Zweiparteienberechnung sowie sichere Mehrparteienberechnung (vgl. Abschnitt 13.2) im Allgemeinen (Beispiele teilweise aus Gentry [34]). Aber auch in unternehmensinternen Big-Data-Systemen könnte homomorphe Verschlüsselung einen erheblichen Schutz von Daten bieten sowie datenschutzkonforme Auswertungen von personenbezogene Daten ermöglichen. Daher gilt homomorphe Verschlüsselung als »der heilige Gral« für Big Data.

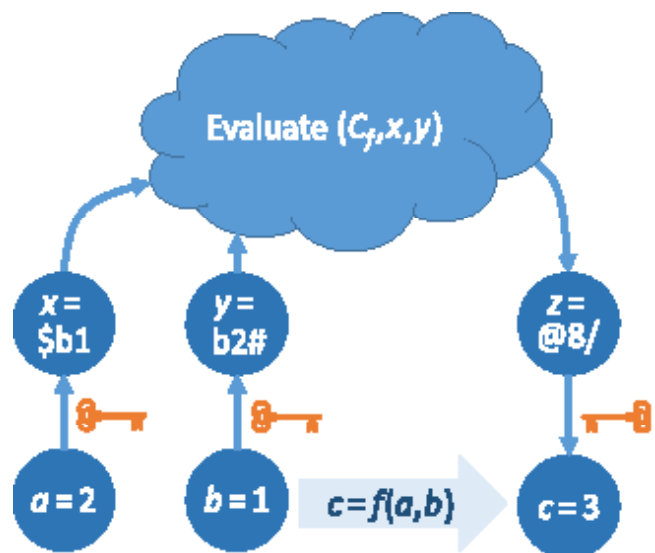


Abbildung 13.1: Homomorphe Verschlüsselung ermöglicht die Datenverarbeitung in der Cloud auf verschlüsselten Inhalten.

Die homomorphe Verschlüsselung ermöglicht die Verarbeitung verschlüsselter Daten aufgrund eines grundlegenden Merkmals: Es gibt Operationen auf den Geheimtexten, die Klartextoperationen entsprechen und einander entsprechende Operationen sind mit der Entschlüsselung vertauschbar, d. h. sie kommutieren mit der Entschlüsselung. Daher erzeugt eine Operation mit verschlüsselten Daten ein verschlüsseltes Ergebnis, das – nach der Entschlüsselung – dem Ergebnis der entsprechenden Operation mit den unverschlüsselten Daten entspricht. Es ist zu beachten, dass die Operation in der verschlüsselten Domäne möglicherweise nicht mit der entsprechenden Operation in der unverschlüsselten Domäne identisch ist. So könnte beispielsweise die Klartextaddition der Geheimtextmultiplikation entsprechen.

Homomorphe Verschlüsselungsverfahren unterscheiden sich in Bezug auf die Art und Anzahl der unterstützten Operationen, die Anzahl der zulässigen aufeinanderfolgenden Operationsauswertungen sowie die Sicherheit und den Aufwand. Eine allgemeine Faustregel ist, dass Verfahren, die mehr Berechnungen ermöglichen, die Sicherheitsherausforderungen und den Rechenaufwand erhöhen.

Die in einem homomorphen Verschlüsselungsschema verfügbaren Operationen können zu komplexeren Berechnungen kombiniert werden. Das Schema kann eine Funktion berechnen, wenn diese Funktion mit einer geeigneten Formel ausgedrückt werden kann, die nur die verfügbaren Operationen verwendet. Aus Computersicht (oder, formaler gesprochen, aus der Perspektive der Berechenbarkeitstheorie und Automatentheorie) entsprechen die Operationen den Gattern, und eine Formel entspricht einer Schaltung. Wenn wir also eine bestimmte Funktion  $f$  auf homomorph verschlüsselten Daten berechnen wollen, müssen wir die Funktion zunächst mit einer geeigneten Formel oder Schaltung  $C_f$  beschreiben. Dann können wir diese Schaltung und die verschlüsselten Eingabedaten einer Prozedur namens Evaluate zur Verfügung stellen, die die Schaltung auf den verschlüsselten Daten mit den entsprechenden Operationen des Geheimtextraums auswertet. Das Ergebnis dieses Verfahrens ist eine verschlüsselte Version des Ergebnisses der Funktion, wie es entstünde, wenn sie auf die unverschlüsselten Eingangsdaten angewendet würde.

$$f(x_1, \dots, x_n) = \text{Decrypt}(\text{Evaluate}(C_f, \text{Encrypt}(x_1), \dots, \text{Encrypt}(x_n))) \quad (13.1)$$

Das höchste Ziel der homomorphen Verschlüsselung ist die Konstruktion von Schemata, die jede Berechnung auf verschlüsselten Daten ermöglichen. Da jede Berechnung mit einer booleschen Schaltung ausgedrückt werden kann, ist das

Ziel, Schemata zu erstellen, die boolesche Schaltungen auswerten können. Dies erfordert die Unterstützung eines geeigneten Satzes von booleschen Operationen im homomorphen Schema. Zum Beispiel reichen die booleschen Operationen AND und NOT aus, um beliebige andere boolesche Operationen auszudrücken und somit eine boolesche Schaltung aufzubauen. Alternativ ist auch die einzige Operation NAND ausreichend.

In Bezug auf das Ziel der Erstellung homomorpher Verschlüsselung sind die Operationen XOR und AND jedoch als Basisoperationen geeigneter, da diese Operationen der Addition und Multiplikation entsprechen, wenn man die zweielementige boolesche Algebra  $\{\text{false}; \text{true}\}$  (Schaltalgebra) als ganze Zahlen  $\{0; 1\}$  mit Durchführung der arithmetischen Operationen Modulo 2 interpretiert.<sup>1</sup> Bei Betrachtung dieser Operationen eröffnet sich eine Vielzahl von mathematischen Objekten als mögliche Geheimtexträume.

Streng genommen reichen XOR und AND nicht aus, um den Aufbau beliebiger boolescher Schaltungen zu ermöglichen. XOR und AND müssen durch die Konstante true (in der Computerperspektive durch ein eingabenloses Gatter, das immer zu true ausgewertet wird) ergänzt werden, um funktional vollständig zu sein. Daher muss ein homomorphes Verschlüsselungsschema neben den Operationen, die XOR (Addition) und AND (Multiplikation) entsprechen, eine kanonische Kodierung der Konstanten true (1) im Geheimtextraum vorsehen.

Ein Schema zu finden, welches beliebige Operationen basierend auf den oben beschriebenen Ideen ermöglicht, war mehrere Jahrzehnte lang der heilige Gral der Kryptografie, bis Gentry [35] 2009 das erste derartige Schema vorschlug (siehe unten).

### 13.1.1 Privatsphären-Homomorphismen

Der erste Ansatz zur Datenverarbeitung auf verschlüsselten Daten wurde bereits 1978 von Rivest et al. [66] eingeführt. Interessanterweise war die Motivation im Jahr 1978 ähnlich wie beim heutigen Cloud Computing: Ein Unternehmen möchte, dass ein Dritter beispielsweise Kreditdaten analysiert, ohne diesem Dritten Einblick in die tatsächlichen Kredite seiner

Kunden zu geben.

Die Grundidee ist ein Homomorphismus, d. h. eine strukturerhaltende Abbildung zwischen zwei algebraischen Strukturen des gleichen Typs, die jeweils aus einer Menge von Objekten, Operationen, Prädikaten (d. h. Relationen wie eine Ordnungsrelation, also ein Größenvergleich) und Konstanten bestehen,

<sup>1</sup> Dieser Ansatz macht die boolesche Algebra zu einer Struktur, die in der Mathematik »Ring« genannt wird. Zhgalkin [88, 89] und Stone [79] waren die ersten, die mit dieser Interpretation boolesche Algebren analysierten. Boyar et al. [12] analysierten die Komplexität boolescher Schaltungen, die mit den Ringoperationen ausgedrückt wurden, in Anbetracht der Bedeutung dieses Themas für die Kryptografie.

wobei jede Operation, jedes Prädikat und jede Konstante ein Eins-zu-Eins-Äquivalent in der anderen Struktur hat. Eine der Strukturen ist der Klartextraum und die andere Struktur ist der Geheimitextraum.

Ein Privatsphären-Homomorphismus ist ein surjektiver Homomorphismus vom Geheimitextraum zum Klartextraum. Somit ist der Privatsphären-Homomorphismus die Entschlüsselungsfunktion. Homomorphismus bedeutet, dass er mit den gegebenen Operationen vertauschbar ist, d. h. mit diesen kommutiert, und die Prädikate und Konstanten bewahrt. Surjektiv zu sein bedeutet, dass die Entschlüsselungsfunktion als Ergebnis jeden beliebigen Klartext liefern kann, was notwendig ist, um die Verschlüsselungsmethode invertieren zu können. Es ist zu beachten, dass die Verschlüsselungsmethode möglicherweise kein Homomorphismus ist (ein solcher inverser Homomorphis-

mus existiert möglicherweise gar nicht), und sie könnte sogar nicht einmal deterministisch sein oder, genauer gesagt, sie sollte aus Sicherheitsgründen nicht deterministisch sein.

Es war von Anfang an klar, dass die Beibehaltung von zu viel Struktur, z. B. einer Ordnungsrelation, es einem Angreifer erlaubt, die Klartexte aus den Geheimitexten offenzulegen. Im Allgemeinen legen Prädikate Informationen offen, da die Auswertung eines Prädikats auf verschlüsselten Daten zu einem unverschlüsselten booleschen Ergebnis führt. Aber auch ohne Prädikate kann ein Angreifer oft die unterstützten Operationen nutzen, um Informationen zu erhalten und schließlich die Verschlüsselung zu brechen. Daher stellte sich lange Zeit die Frage, ob es Arten von algebraischen Strukturen gibt, die einerseits für Berechnungen nützlich genug sind und andererseits sichere Privatsphären-Homomorphismen ermöglichen.

### 13.1.2 Partiiell homomorphe Verschlüsselung (PHE)

Wenn ein Verschlüsselungsschema nur mit einer von mehreren Operationen, z. B. Addition oder Multiplikation, kommutiert, wird es partiell homomorphe Verschlüsselung (PHE) genannt. Beispiele sind RSA für Multiplikationen [67], das Goldwasser-Micali-Kryptosystem für XOR-Operationen [36] oder das Benaloh-Kryptosystem für Additionen und Subtraktionen [7]. Aufgrund fehlender Unterstützung für einige der gewünschten Operationen waren solche Verschlüsselungsschemata nicht für nicht-interaktive Schaltungsauswertungen geeignet, die für Cloud Computing auf verschlüsselten Daten erforderlich wären, aber sie wurden für die Verwendung in interaktiven Protokollen für Zero-Knowledge-Beweise (z. B. [13]) oder in der interaktiven Auswertung von kryptografisch »verwüfelten« Schaltungen (engl. garbled circuits, z. B. [1]) vorgeschlagen.

Stephen, Savvides et al. haben Crypsis [78] und den Nachfolger Cuttlefish [71] eingeführt, zwei Prototypen, die unter anderem PHE nutzt, um eine Big-Data-Verarbeitung direkt auf verschlüsselten Daten zu ermöglichen. Crypsis ist für Pig-Latin-Skripte konzipiert, die auf einer MapReduce-Infrastruktur ausgeführt werden. Crypsis analysiert zunächst in einem

solchen Skript, welche Operatoren auf welche Daten und Konstanten angewendet werden. Dann wendet es ein entsprechendes Verschlüsselungsschema auf die entsprechenden Daten an und transformiert das Skript entsprechend. Cuttlefish verarbeitet Routinen, die in Spark geschrieben sind mit einem eigenen Compiler und einer Planner Engine, um bei der Ausführung eine möglichst hohe Effizienz zu erreichen. Beide Systeme unterstützen additiv homomorphe Verschlüsselung, multiplikativ homomorphe Verschlüsselung, deterministische Verschlüsselung zur Unterstützung des Gleichheitsprädikats und ordnungserhaltende Verschlüsselung für das Ordnungsprädikat (Größenvergleich). Bei Bedarf werden die Daten mit mehreren Schemata verschlüsselt oder es werden Wiederverschlüsselungsroutinen in das Skript eingefügt. Die Verschlüsselung, Entschlüsselung und Wiederverschlüsselung erfolgt in einer vertrauenswürdigen Umgebung, während das transformierte Skript in einer unsicheren Umgebung ausgeführt werden kann. Die Auswirkungen auf die Sicherheit durch Mehrfachverschlüsselung derselben Daten und durch die Wiederverschlüsselung müssen noch analysiert werden.

### 13.1.3 »Begrenzt« homomorphe Verschlüsselung (SHE)

Begrenzt homomorphe Verschlüsselung (engl. somewhat homomorphic encryption, SHE) erlaubt im Gegensatz zu PHE sowohl Addition als auch Multiplikation, aber jede Operation verursacht einen Rechenfehler oder ein »Rauschen«, das sich mit jeder weiteren Operation akkumuliert, bis das Ergebnis nicht mehr eindeutig entschlüsselt werden kann. Dies begrenzt die Anzahl der aufeinanderfolgenden Operationen, die mit verschlüsselten Daten durchgeführt werden können, oder, unter Verwendung der Computerterminologie, die Tiefe der auszuwertenden Schaltungen. Das Schema von Sander et al. [70] ermöglicht die Auswertung von Schaltungen mit logarithmischer Tiefe. Das von Boneh et al. [11] vorgeschlagene Schema erlaubt die Auswertung einer beliebigen Anzahl von Additionen, aber höchstens einer Multiplikation.

#### Voll-homomorphe Verschlüsselung (FHE)

Wenn ein homomorphes Verschlüsselungsschema die Addition und Multiplikation auf verschlüsselten Daten ohne die Mängel von SHE unterstützt, d. h. ohne die Tiefe der auszuwertenden Schaltungen einzuschränken, spricht man von einer voll-homomorphen Verschlüsselung (engl. fully homomorphic encryption, FHE). Das erste FHE-Schema wurde von Gentry im Jahr 2009 veröffentlicht [34, 35].

Die Grundidee von Gentrys Konstruktion ist eine Strategie, die er Bootstrapping nannte. In einem ersten Schritt konstruierte er ein geeignetes SHE-Schema. Bootstrapping verwandelt ein SHE-Schema in ein FHE-Schema, indem es Daten neu verschlüsselt, bevor das durch die Operationen des SHE-Schemas verursachte Rauschen zu groß wird. Dazu müssen die verschlüsselten Daten noch einmal verschlüsselt werden und die Entschlüsselungsschaltung muss auf diesen doppelt verschlüsselten Daten ausgewertet werden. Diese Schaltung selbst muss einfach genug sein, um nicht zu viel Rauschen einzubringen, bevor sie fertig ist. Dann ergibt sich ein neuer Geheimtext, der den ursprünglich verschlüsselten Daten entspricht, aber mit so wenig Rauschen wie frisch verschlüsselte Daten. Leider war die Entschlüsselungsschaltung von Gentrys SHE-System etwas zu komplex. Daher führte Gentry zusätzlich eine Strategie ein, die er Squashing nannte, um die Entschlüsselungsschaltung auf Kosten eines komplexeren Verschlüsselungsprozesses zu vereinfachen. Dies machte sein SHE-Schema Bootstrapping-fähig. Nun kann jedes Gatter in einer beliebigen Schaltung von einem Wiederverschlüsselungsschritt

basierend auf der »gesquashten« Entschlüsselungsschaltung begleitet werden, um sicherzustellen, dass das Rauschen nie zu groß wird.

Kurz nach Gentrys erstem FHE-Schema stellten van Dijk et al. [20] ein weiteres FHE-Schema mit einer einfacheren SHE in seinem Inneren vor. Der Rest ihres Schemas verwendet die gleichen Ideen wie das ursprüngliche Schema von Gentry, nämlich Squashing und Bootstrapping. Das Schema von van Dijk et al. kann nun als kanonisches Beispiel für FHE betrachtet werden.

Nach der Definition von Gentry ist FHE eine Public-Key-Verschlüsselung. Das bedeutet, dass sich der für die Verschlüsselung verwendete Schlüssel von dem für die Entschlüsselung verwendeten Schlüssel unterscheidet, und nur der Entschlüsselungsschlüssel geheim gehalten werden muss, während der Verschlüsselungsschlüssel veröffentlicht werden kann. Somit kann jeder verschlüsseln, aber nur der Besitzer des privaten Schlüssels kann entschlüsseln. Dies ermöglicht mehr Nutzungsszenarien als symmetrische Verschlüsselungsverfahren, bei denen derselbe Schlüssel zur Ver- und Entschlüsselung verwendet wird. So ermöglicht Public Key FHE beispielsweise Szenarien, in denen auch Fremdeingaben in die Verarbeitung verschlüsselter Daten einbezogen werden können. Das von van Dijk et al. [20] vorgeschlagene Schema hat eine symmetrische Variante und eine Public-Key-Variante. Die meisten PHE-Schemata (vgl. oben) sind ebenfalls Public-Key-Schemata.

Alle bestehenden FHE-Systeme haben den Nachteil, dass sie komplex in der Konstruktion sind und einen erheblichen Mehraufwand für die Berechnung bis zu einem für die Praxis nicht akzeptablen Niveau verursachen. Während FHE die wünschenswerteste Variante der homomorphen Verschlüsselung ist, gelten nur PHE und SHE als in der Praxis verwendbar. Naehrig et al. [43] diskutieren die Kosten für die Nutzung eines SHE-Systems im Detail.

Man beachte, dass es einen wichtigen Unterschied zu Privatsphären-Homomorphismen bei der Prädikatsauswertung gibt. Sowohl ein Privatsphären-Homomorphismus als auch ein FHE-Schema können Prädikate auswerten, aber der Unterschied liegt in der Form, in der das Ergebnis bereitgestellt wird. Wenn ein Privatsphären-Homomorphismus ein Prädikat unterstützt, führt die Auswertung dieses Prädikats auf verschlüsselten

Daten zu einer Klartextantwort. Ein FHE-Schema kann jedes Prädikat auswerten, wenn eine Schaltung zur Verfügung gestellt wird, die dieses Prädikat implementiert.<sup>2</sup> Das Ergebnis

dieser Prädikatsberechnung liegt jedoch in der verschlüsselten Domäne. Somit gibt es einerseits kein Informationsleck, andererseits kann das Ergebnis aber nicht direkt genutzt werden.

## 13.2 Sichere Mehrparteienberechnung

Sichere Mehrparteienberechnung (engl. secure multiparty computation, MPC) kann als das gemeinschaftliche Berechnen einer Funktion verstanden werden, für die mehrere Parteien private Eingaben liefern, um ein Ergebnis über ihre gesammelten Eingaben zu berechnen. Dabei werden die individuellen Informationen der jeweiligen Parteien nicht an die übrigen Parteien preisgegeben [47]. Das bedeutet, dass jede Partei – und idealerweise jeder Gegner – nicht mehr aus der gemeinsamen Berechnung lernen kann als aus der eigenen Eingabe und dem Endergebnis, d. h. der Funktionsausgabe.

Ein beispielhaftes Szenario für die Anwendung sicherer Mehrparteienberechnungen ist die gemeinsame Nutzung medizinischer Datenbanken. Man stelle sich mehrere Patientendatenbanken vor, die jeweils private und sensible Informationen über Krankheitsdiagnosen enthalten. Das Zusammentragen von Informationen über die Vereinigungsmenge der Datenbanken könnte wichtige Erkenntnisse für Forschung und das Gesundheitssystem liefern. Allerdings würde man durch das Teilen der Datenbanken private Informationen der Patienten preisgeben, was jedoch nicht ohne Weiteres zulässig ist.

Ein anderes denkbare Szenario sind Wahlen, bei denen man wissen will, welcher Kandidat letztlich gewählt wurde, aber nicht, wer für welchen Kandidaten gestimmt hat. Ein klassisches, berühmtes Beispiel für sichere Mehrparteienberechnung ist Yaos Millionärsproblem [87]: Zwei (oder mehr) Millionäre möchten eruieren, wer der reichste von ihnen ist, aber keiner möchte sagen, wie viel Geld er tatsächlich besitzt. Ein recht junges, aber gerade im Kontext von Big Data relevantes Szenario ist maschinelles Lernen auf den Daten vieler Personen. Hierzu sind in den letzten Jahren Vorschläge entstanden, die Privatsphäre der Betroffenen mithilfe von MPC zu schützen, s. Kapitel 16.

Eine typische und technisch ausführlich erforschte Aufgabe im Kontext von MPC ist die Bestimmung von Schnittmengen.

In der Praxis kann es sich bei den Mengen beispielsweise um Personenlisten handeln. So kann man etwa die Passagierlisten von Fluggesellschaften mit Fahndungslisten von Behörden abgleichen, ohne die Passagierlisten an die Behörden oder die Fahndungslisten an die Fluggesellschaften geben zu müssen. Auf diese Weise wird ein datenschutzfreundlicher Abgleich ermöglicht. Die beteiligten Parteien erfahren mit MPC ausschließlich, welche Personen in der Schnittmenge enthalten sind. Konkret erfährt die Fluggesellschaft, nach welchen Passagieren gefahndet wird und die Behörden erfahren, welche gesuchte Person einen Flug gebucht hat. Aber die Behörden erfahren nicht über jede beliebige Person, welchen Flug sie gebucht hat und die Fluggesellschaften erfahren nicht, welche Personen insgesamt auf den Fahndungslisten stehen.

Das Ziel eines sicheren Mehrparteienberechnungsprotokolls ist, dass jede beteiligte Partei nicht mehr weiß als das, was sie selbst als Eingabe geliefert hat, d. h. ihre eigenen Datenbankinformationen, die zur Berechnung beigetragen haben, und das finale Ergebnis.

Um die Sicherheit eines MPC-Protokolls zu beweisen, wird gedanklich ein Vergleich mit einem sogenannten idealen Modell durchgeführt. Im idealen Modell senden die beteiligten Parteien ihre Eingaben auf einem sicheren Weg an eine vertrauenswürdige Instanz, die die Berechnung durchführt und die Ausgaben anschließend an die Parteien zurück sendet. Im realen Modell erhalten Parteien keine Hilfe von einer vertrauenswürdigen Instanz, sondern sie führen das MPC-Protokoll aus. Ein Protokoll wird als sicher bezeichnet, wenn jeder Angriff, den man sich im realen Modell erdenken kann, ebenso im Idealmodell durchgeführt werden kann. Es können jedoch keine Angriffe auf das ideale Modell ausgeführt werden, weshalb Sicherheit direkt impliziert wird. Somit kann ein Angreifer beim Abhören oder Sabotieren einer Mehrparteienberechnung höchstens das herausfinden, was er auch direkt aus dem finalen Ergebnis und den Eingaben der von ihm kontrollierten Parteien lernen könnte.

<sup>2</sup> Dies ist in der Tat nicht auf FHE beschränkt. Weniger leistungsfähige Schemata könnten auch in der Lage sein, das Prädikat auszuwerten, wenn das Prädikat mit einer nach diesem Schema auswertbaren Schaltung ausgedrückt werden kann. Aber mit FHE wissen wir, dass jede Schaltung ausgewertet werden kann, und somit jedes entscheidbare (berechenbare) Prädikat ausgewertet werden kann.

Formal wird die Definition der rechnerischen Nichtunterscheidbarkeit verwendet:

**Definition 13.2.1 (Rechnerische Nichtunterscheidbarkeit)**

*Ein MPC-Protokoll ist sicher, wenn eine Simulation der Parteien (inkl. Angreifer) und ihrer Kommunikation auf Basis der Informationen aus dem ideale Modell von einer Ausführung des Protokolls im realen Modell rechnerisch nicht unterscheidbar ist. Die Simulation und die reale Ausführung werden dabei als rechnerisch nicht unterscheidbar bezeichnet, wenn jeder probabilistische (d. h. nichtdeterministische) Polynomialzeit-Beobachter in beiden Fällen bis auf eine vernachlässigbare Abweichung (die mit wachsendem Sicherheitsparameter gegen Null strebt) mit gleicher Wahrscheinlichkeit eine 1 ausgibt.*

Mathematisch tiefergehende Ausarbeitungen hiervon finden sich bei Lindell und Pinkas [48].

Erläuterung: Ein Beobachter kennt alle Ein- und Ausgaben, insbesondere auch das Wissen, das ein Angreifer aus der Kontrolle von Parteien oder dem Abhören von Kommunikation berechnet. Der Beobachter ist hier ein „Supercomputer“, der alle Probleme in der Komplexitätsklasse NP effizient lösen kann (bzw. alle NP-Algorithmen effizient ausführen kann). Ein solcher Beobachter hat bei einem sicheren Protokoll auch mit dem am besten geeigneten Algorithmus innerhalb dieser Klasse nur eine vernachlässigbare Chance, eine Simulation

des Protokolls auf Basis des idealen Modells von einer echten Ausführung des Protokolls im realen Modell zu unterscheiden. Dabei wird ohne Beschränkung der Allgemeinheit angenommen, dass der Beobachter eine 1 ausgibt, wenn er vermutet, dass es sich um eine Simulation handelt und eine 0, wenn er vermutet, dass es sich um eine reale Ausführung handelt. Wenn diese Unterscheidung nicht gelingt, sind die Erkenntnisse des Angreifers im realen Modell nicht wesentlich anders als die Erkenntnisse des simulierten Angreifers im idealen Modell. Daher kann der Angreifer im realen Modell nicht mehr lernen als die möglichen Erkenntnisse aus dem, was ihm im idealen Modell zur Verfügung steht, also aus den Eingaben der von ihm kontrollierten Parteien und aus dem Endergebnis.

Im Vergleich zu homomorpher Verschlüsselung ist MPC effizienter und in einigen Szenarien bereits praxistauglich. Der Nachteil von MPC ist allerdings, dass für jede Art von Berechnung, etwa die Summenberechnung oder die Schnittmengenbestimmung, ein eigenes MPC-Protokoll entworfen werden muss, während voll-homomorphe Verschlüsselung theoretisch universell für alle Berechnungen einsetzbar ist. Ansätze, die MPC generischer machen, indem die Auswertung verschiedener Schaltungen zugelassen wird, nutzen oft homomorphe Verschlüsselung als Hilfsmittel und sind dementsprechend ineffizient. Auch andere generische Ansätze für MPC sind recht ineffizient.

## 14 Anonymisierung strukturierter Daten

Für eine Datenverarbeitung sind strukturierte Daten von Vorteil. Daher werden Daten, wenn möglich, in strukturierter Form der Verarbeitung zugeführt. Strukturierte Daten können tabellarisch dargestellt werden, wobei meist jede Spalte ein Attribut enthält und jede Zeile einen Datensatz. Bei personenbezogenen Tabellen ist je eine Zeile einer Person zugeordnet.

Das Ziel der Wahrung der Privatsphäre wird erreicht, indem sensible Daten oder unnötig viele Informationen, die zum Verknüpfen mehrerer Datenbestände genutzt werden könnten, nicht offengelegt werden, d. h. es werden nicht mehr Inhalte an andere Parteien weitergegeben als erforderlich. Hierbei müssen verschiedene Begriffe differenziert werden, konkret Privatheit und Anonymität. Informationen sind dann privat, wenn sie nicht von Personen, die nicht der entsprechenden Privat-»sphäre« angehören, eingesehen werden können. Aus technologischer Sicht bedeutet dies, dass die Informationen vertraulich sind. Anonymität hingegen verhindert, dass Informationen – auch solche, die öffentlich zugänglich sind – mit einer bestimmten Person verknüpft werden können. Dies dient ebenso dem Schutz der Privatsphäre, da die Informationen über eine konkrete Person nicht abgerufen werden können. Dabei muss man jedoch beachten, dass es auch dann Risiken für die Privatsphäre gibt, wenn der Kreis in Frage kommender Personen stark eingeschränkt werden kann (s. Unterabschnitt 17.2.2 für eine weitergehende Diskussion dieses Aspekts). Daher wird in der Wissenschaft Anonymität eher als Maß für den Grad der Geheimhaltung der Identität einer Person angesehen.

Für die Anonymisierung strukturierter personenbezogener Daten gibt es verschiedene elementare Strategien:

### Generalisierung:

Die jeweiligen Attributwerte werden durch weniger genaue Angaben ersetzt, etwa durch Intervalle bei numerischen Daten oder durch übergeordnete Kategorien bei kategorischen Daten.

### Löschung:

Der Inhalt einzelner Zellen, Spalten oder Zeilen wird gelöscht. Dies entspricht einer Generalisierung zu einem allumfassenden und nichtssagenden Wert, etwa „\*“.

### Mikroaggregation:

Die Daten werden nach Ähnlichkeit in den Attributwerten gruppiert (engl. clustering), und pro Gruppe werden die einzelnen Werte zu einem repräsentativen Wert zusammengefasst, etwa dem Mittelwert oder Median.

### Verfälschung:

Ein Teil der Daten oder alle Daten werden zufällig abgewandelt. Dies kann z. B. dadurch erreicht werden, dass zu den Werten zufällige Störungen hinzugefügt werden, dass verschiedene Einträge in der Tabelle vertauscht werden oder dass eine künstliche Tabelle unter Orientierung an der Originaltabelle synthetisiert wird.

Neben diesen verschiedenen Ansätzen zur Anonymisierung der Daten gibt es auch die Strategie, nicht die Daten in anonymisierter Form herauszugeben, sondern die gewünschte Analyse zu den Originaldaten in geschützter Umgebung zu bringen, die Analyse dort durchzuführen und nur die Ergebnisse vor der Herausgabe zu anonymisieren. Dies kann einfacher sein und präzisere Ergebnisse liefern, aber es muss stets die rechtliche Zulässigkeit einer solchen Verarbeitung geprüft werden. Zudem schwindet der Vorteil und kann sich in das Gegenteil kehren, wenn viele Analysen durchgeführt werden sollen, da bei der Anonymisierung der Ergebnisse dann auch die Querbeziehungen zwischen allen Ergebnissen berücksichtigt werden müssen.

Zum Bestimmen des Anonymitätsgrades von Daten, die mit den oben genannten Strategien behandelt worden sind, gibt es verschiedene Kriterien bzw. Maße. Diese Maße unterscheiden sich darin, welche Annahmen über das Hintergrundwissen eines Angreifers und über die Art des zu erreichenden Schutzes gemacht werden. Minimalen Schutz bieten die Kriterien  $k$ -Map [82] und  $\delta$ -Presence [59], da hier angenommen wird, dass die in der Tabelle erfassten Individuen aus einer größeren Population stammen, ein Angreifer aber nicht wissen kann, ob eine bestimmte Person in der Tabelle enthalten ist. Das bekannteste Anonymitätskriterium ist  $k$ -Anonymität [82]. Nach diesem Kriterium muss es jeweils mindestens  $k$  für eine Person in Frage kommende Einträge in der Tabelle geben, sodass eine Re-Identifikation nicht möglich ist. Da dennoch möglicherweise einzelne Attribute einer Person durch eine



k-anonyme Tabelle offengelegt werden können, wurde das Kriterium zu I-Diversität [51] und t-Nähe (engl. t-closeness) [46] weiterentwickelt. k-Anonymität, I-Diversität und t-Nähe werden in Abschnitt 14.1 genauer betrachtet. Grundsätzlich anders ist das Konzept von Differential Privacy [27]. Hier wird

die Anonymität daran gemessen, wie sehr sich das Ergebnis durch Weglassen oder Hinzufügen einer Person ändern kann, und somit wie viel an Information maximal über eine Person offenbart wird. Dieses Konzept wird in Abschnitt 14.2 ausführlicher erläutert und analysiert.

## 14.1 Anonymisierung auf Basis von k-Anonymität

Im Folgenden erläutert dieser Bericht die Konzepte von k-Anonymität, I-Diversität und t-Nähe (t-closeness). Diese Konzepte haben ein gemeinsames Fundament. Sie bauen auf der Sichtweise auf, dass die in Datensätzen enthaltenen Attribute in folgende drei Klassen eingeteilt werden können [52]:

### Explizite Identifikatoren:

Explizite Identifikatoren sind Attribute des Datenbestands, die die entsprechende Person sofort identifizieren. Beispiele sind Vor- und Nachname, Ausweisnummer oder Sozialversicherungsnummer.

### Quasi-Identifikatoren:

Quasi-Identifikatoren sind Mengen von Attributen, die mit geringem Aufwand oder durch Verknüpfungen die entsprechende Person möglicherweise identifizieren könnten. Ob eine Menge von Attributen ein Quasi-Identifikator ist, hängt vom gegebenen Szenario ab. Beispiele für Quasi-Identifikatoren in der Datenbank eines Unternehmens sind Geburtstag, Einkom-

men oder Tag des Eintritts in das Unternehmen.

### Sensible Attribute:

Ein sensibles Attribut beschreibt ein personenspezifisches Merkmal, welches bei Offenlegung die Privatsphäre der jeweiligen Person verletzt. Ein sensibles Attribut kann beispielsweise Krankheit, Behinderung, Einkommen, Einkaufshistorie oder Kreditwürdigkeit sein.

Die nachfolgend vorgestellten Anonymitätsmaße können nur dann die Privatsphäre der in einer Tabelle enthaltenen Personen schützen, wenn die Zuordnung der Attribute zu diesen drei Kategorien korrekt ist, d. h. konsistent mit den Möglichkeiten eines Angreifers in der Praxis. Sonst bleibt eine Re-Identifizierung möglich und die Anonymisierung ist hinfällig. Daher ist ein gutes Wissen über Art, Aufbau und Inhalt der Daten erforderlich, um die Anonymisierungstechniken auf geeignete Weise auf große Datensätze anwenden zu können.

---

### 14.1.1 k-Anonymität

---

k-Anonymität ist ein Kriterium zur *Privatheit-erhaltenden Datenveröffentlichung* (engl. privacy-preserving data publishing (PPDP)), welches 1998 von Samarati und Sweeney erstmals vorgestellt [69] und seit Sweeney's Veröffentlichungen im Jahr 2002 [80, 81] sehr populär wurde. Sweeney demonstrierte, wie die Krankenakte einer bestimmten Person (William Weld, Gouverneur von Massachusetts) auf Grundlage von Quasi-Identifikatoren eindeutig identifiziert werden konnte, indem sie eine angeblich anonymisierte (d. h. alle expliziten Identifikatoren wurden gelöscht) medizinische Datenbank mit einer Wählerregistrierungsliste verknüpfte [80].

#### Grundsätze

k-Anonymität anonymisiert Daten, indem eine Person von  $k - 1$  weiteren ununterscheidbar gemacht wird. Demnach

sprechen wir von k-Anonymität, wenn für jede Person mindestens k Datensätze in Frage kommen. Dementsprechend gilt: Je größer k, desto höher ist der Grad der Anonymität der Gruppe. Jede Gruppe von mindestens k Personen, die den selben Quasi-Identifikator aufweisen, wird als Äquivalenzklasse bezeichnet.

#### Angriffe und Limitationen

k-Anonymität kann keinen 100-prozentigen Schutz der Privatheit gewährleisten [52]. Bestimmte Angriffe, wie beispielsweise ein Homogenitätsangriff [84] oder eine Background Knowledge Attack [52], können private Informationen, d.h. sensible Attribute einer bekannten einzelnen Person oder einer Gruppe, enthüllen. Ein weiterer Kritikpunkt an Algorithmen für k-Anonymität ist die Komplexität der Berechnung,

insbesondere in Fällen, in denen die Anzahl der Attribute und die Gesamtgröße der Datenbank hoch ist (vgl. Unterabschnitt

14.1.4).

### 14.1.2 I-Diversität

2006 stellten Machanavajjhala et al. [51, 52] Angriffe auf das Modell der  $k$ -Anonymität vor, welche auf Homogenität bzw. auf Hintergrundwissen basierten. Um diese Schwachstellen von  $k$ -Anonymität zu beheben, führten sie das Konzept der I-Diversität ein.

#### Grundsätze

I-Diversität greift an den Werten der sensiblen Attribute an. In einer Äquivalenzklasse, in der die Werte eines sensiblen Attributs alle gleich oder unzureichend indifferent sind, ist dieses Attribut nicht mehr davor sicher, offengelegt zu werden. Daher ist eine Äquivalenzklasse I-divers, wenn alle empfindlichen Attribute mindestens  $l$  »stark vertretene« Werte enthalten [52]. Eine Datenbank ist dann I-divers, wenn alle ihre Äquivalenzklassen I-divers sind. Es gibt verschiedene Auffassungen darüber, was genau »stark vertreten« in diesem Kontext bedeutet.

#### Eindeutige I-Diversität (engl. distinct l-diversity):

Hier bezieht sich »stark vertreten« auf die Anzahl der verschiedenen Werte. Eine Tabelle ist folglich eindeutig I-divers, wenn es mindestens  $l$  verschiedene Werte für jedes sensible Attribut gibt. Unterscheiden sich diese Werte für jedes sensible Attribut stark voneinander, so gilt dieser Ansatz als schwach.

$$H(EC) = - \sum_{s \in S} \mathbb{P}(s|EC) \log \mathbb{P}(s|EC),$$

#### Entropie-I-Diversität:

Eine Interpretation von I-Diversität, welche starken Schutz bietet, ist die Entropie-I-Diversität. Sei  $EC$  eine Äquivalenzklasse, die Entropie von  $EC$  wird dann wie folgt definiert:

wobei  $S$  alle Werte repräsentiert, die das sensible Attribut annehmen kann, und  $\mathbb{P}(s|EC)$  die relative Häufigkeit der Einträge in  $EC$  darstellt, deren sensibles Attribut die Ausprä-

gung  $s$  annimmt. Eine Datenbank ist demnach Entropie-I-divers, wenn  $H(EC) \geq \log l$  für alle Äquivalenzklassen  $EC$  in der Datenbank gilt.

#### Rekursive (c, l)-Diversität:

Ein Kompromiss zwischen Eindeutiger I-Diversität und Entropie-I-Diversität wird durch die sogenannte rekursive  $(c, l)$ -Diversität gebildet. Die Idee hinter diesem Ansatz ist sicherzustellen, dass der häufigste Wert in einem sensiblen Attribut ausreichend selten und gleichzeitig der seltenste Wert ausreichend häufig vertreten ist. Daher: Seien  $s_1, \dots, s_m$  mögliche Ausprägungen für ein sensibles Attribut und  $r_i$  für  $1 \leq i \leq m$  die zugehörigen Ausprägungshäufigkeiten innerhalb einer Äquivalenzklasse, geordnet von der häufigsten bis zur seltensten Ausprägung. Gegeben eine Konstanten  $c$  gilt eine Äquivalenzklasse als rekursiv  $(c, l)$ -divers, wenn  $r_1 < c(r_1 + r_1 + 1 + \dots + r_m)$ . Eine Datenbank gilt dann als rekursiv  $(c, l)$ -divers, wenn alle ihre Äquivalenzklassen rekursiv  $(c, l)$ -divers sind.

#### Angriffe und Limitationen

Wenn ein sensibles Attribut dadurch gekennzeichnet ist, dass es eine geringe Anzahl an Ausreißern aufweist, erhöht I-Diversität die Privatheit nicht unbedingt. Wenn beispielsweise nur 1 von 1000 Personen in einer Tabelle HIV-positiv ist, aber in einer 2-diversen Äquivalenzklasse die Hälfte der Personen HIV-positiv ist, dann bedeutet das eine signifikante Stigmatisierung aller Personen in dieser Äquivalenzklasse. Ist hingegen ein sensibles Attribut über den gesamten Datensatz hinweg konstant, so ist I-Diversität überflüssig.

Das Kriterium der I-Diversität berücksichtigt auch nicht die semantischen Bedeutungen sensibler Werte, was es ermöglicht, Ähnlichkeitsangriffe durchzuführen. Bei einem Ähnlichkeitsangriff kann der Angreifer ein sensibles Attribut der Zielperson auf wenige, ähnliche Werte einschränken.

### 14.1.3 t-Nähe

Die t-Nähe (engl. t-closeness) versucht, die Schwächen von k-Anonymität und l-Diversität zu überwinden, indem sie zu jeder Äquivalenzklasse ein zusätzliches Kriterium hinzufügt. Genauer gesagt wird das Kriterium von k-Anonymität durch ein zusätzliches Kriterium ergänzt, welches das Kriterium von l-Diversität ersetzt.

#### Grundsätze

Das von Li et al. [46] eingeführte Maß der t-Nähe definiert Privatheit als den Informationsgewinn, den ein Beobachter durch das Ausführen einer Aktion auf der k-anonymisierten Datenbank realisieren kann. Informationsgewinn wird im Allgemeinen als die Differenz zwischen dem Wissen über eine Person vor (a priori) und nach (a posteriori) einem Angriff beschrieben. Um den Informationsgewinn so gering wie möglich zu halten, sollte die Verteilung der sensitiven Attribute innerhalb der jeweiligen Äquivalenzklassen jener über die gesamte Datenbank hinweg so weit wie möglich ähneln.

Li et al. definieren t-Nähe wie folgt: »An equivalence class is said to have t-closeness if the distance between the distribu-

tion of a sensitive attribute in this class and the distribution of the attribute in the whole table is no more than a threshold  $t$ . A table is said to have t-closeness if all equivalence classes have t-closeness.« [46].

#### Angriffe und Limitationen

Angriffe, die spezifisch für t-Nähe sind, sind nicht bekannt. Allerdings besteht bei t-Nähe wie auch bei k-Anonymität und l-Diversität das Problem, dass die Zuordnung der Attribute zu Identifikatoren, Quasi-Identifikatoren und sensiblen Attributen korrekt sein muss, da sonst die Anonymisierung hinfällig ist, s. oben.

Ein anderes Problem von t-Nähe ist, dass ein sehr großer Datenverlust erzwungen wird, da jede Äquivalenzklasse in ihren statistischen Eigenschaften nicht wesentlich von der Grundgesamtheit abweichen darf. Daher werden die Daten schnell wertlos für jegliche Analysen oder Erkenntnisse. Nur wenn der Parameter  $t$  in Abhängigkeit der Domäne und der Daten mit Bedacht gewählt wird, ist eine brauchbare Restinformation bei gleichzeitigem Schutz der Privatsphäre möglich.

### 14.1.4 Algorithmen für k-Anonymität und verwandte Kriterien

Für k-Anonymität und die darauf aufbauenden Kriterien gibt es eine Vielzahl von Algorithmen. Diese Algorithmen entfernen die identifizierenden Attribute (explizite Identifikatoren), während die sensiblen Attribute beibehalten werden. Die wichtigste Transformation passiert mit den Quasi-Identifikatoren, und die existierenden Algorithmen unterscheiden sich hinsichtlich der Strategien zur Anonymisierung der Quasi-Identifikatoren sowie bzgl. des Verfahrens zur Suche nach einer guten Umsetzung der gewählten Strategie.

Hier werden oft Verfahren der Generalisierung und Löschung angewendet, um sicherzustellen, dass jedes Individuum in Bezug auf den Quasi-Identifikator identisch mit mindestens  $k - 1$  anderen ist. Die einfacheren Algorithmen beschränken sich auf eine Generalisierung auf der Attribut-Ebene, d. h. es wird für eine Tabellenspalte insgesamt festgelegt, welcher Wert zu welchem generalisiert wird („global recoding“), während komplexere Algorithmen die Generalisierung auf Zell-Ebene festlegen können („local recoding“). Die zweite Gruppe von Algorithmen kann das Ziel mit weniger Informationsverlust er-

reichen, jedoch ist die Durchführung bei realen Tabellen meist zu aufwendig, da der Aufwand zum Finden einer optimalen Generalisierung bei naiver Suche exponentiell mit der Anzahl der Tabellenzellen steigt. Diese Optimierungsaufgabe ist in der Tat NP-schwer [24], sodass keine effizienten Algorithmen existieren. Aber selbst die erste Gruppe von Algorithmen kann bei großen Datentabellen, insbesondere, wenn viele Attribute vorhanden sind, zu aufwendig werden.

Eine weitere Gruppe von Algorithmen verwendet Mikroaggregation. Algorithmen auf Basis von Mikroaggregation können eine relativ gute Effizienz aufweisen und gleichzeitig mehr Informationen erhalten als Generalisierungen auf Attribut-Ebene. Die wesentliche Herausforderung bei der Mikroaggregation ist das Auffinden optimaler Cluster, d. h. Gruppen von Individuen, die zusammengefasst werden sollen. Etabliert ist hier insbesondere die Familie von MDAV-Algorithmen. MDAV steht für Maximum Distance to Average Vector. Das ursprüngliche Verfahren wurde von Domingo-Ferrer und Mateo-Sanz entworfen [22]. Eine Verbesserung gegenüber

bestehenden MDAV-Varianten wird derzeit von Antonio Odoguardi in seiner Masterarbeit [60] an der TU Darmstadt in Kooperation mit dem Fraunhofer SIT entwickelt.

Ein Nachteil von Mikroaggregation ist, dass diese Strategie nur gut auf Attribute mit kontinuierlichem Wertebereich anwendbar ist, aber schlecht auf kategoriale Attribute. Daher

hat Jamal Pasha in seiner Masterarbeit [62] an der TU Darmstadt in Kooperation mit dem Fraunhofer SIT das Verfahren ( $\alpha, \gamma$ )-Anonymisierung entworfen, welches Mikroaggregation auf kontinuierlichen Attributen mit Generalisierung auf kategorischen Attributen verknüpft. Dieses Verfahren wurde im Rahmen des Projekts »Privacy und Big Data« wissenschaftlich publiziert [72].

## 14.2 Differential Privacy

Das Konzept von *Differential Privacy* wurde von Dwork et al. [29] im Jahr 2006 unter dem Begriff *Ununterscheidbarkeit* (engl. indistinguishability) eingeführt. Ein häufig betrachtetes Szenario besteht darin, dass ein vertrauenswürdiger Datenverarbeiter eine Datenbank mit sensiblen Informationen besitzt und statistische Anfragen an die Datenbank beantwortet. Das Ziel von Differential Privacy ist es, die Menge an Informationen, die über einzelne Datenbankeinträge, d. h. in der Regel über einzelne Personen, aus statistischen Auswertungen der Datenbank gelernt werden können, zu begrenzen.

Differential Privacy ist keine fixe Methode, sondern vielmehr eine Eigenschaft, welche verlangt, dass es für eine beliebige Analyse irrelevant ist, ob ein bestimmtes Datensubjekt im Datensatz enthalten ist oder nicht – in beiden Fällen sollten sich die Output-Verteilungen nicht signifikant unterscheiden. Intuitiv bedeutet dies, dass die Menge an Informationen, die über ein bestimmtes Individuum herausgefunden werden kann, limitiert wird.

Dwork charakterisiert Differential Privacy wie folgt:

»Roughly speaking, differential privacy ensures that the removal or addition of a single database item does not (substantially) affect the outcome of any analysis. It follows that no risk is incurred by joining the database, providing a mathematically rigorous means of coping with the fact that distributional information may be disclosive.« [28]

Differential Privacy wird in der Regel durch das Hinzufügen von Rauschen erreicht. Der Grad der Störung hängt von der Stärke des Einflusses des einzelnen Eintrags auf den Datensatz ab. Um den Beitrag einzelner Einträge zu verbergen, muss das Rauschen so gewählt werden, dass Datenbanken, die sich nur um eine Person unterscheiden, sehr ähnliche Verteilungen möglicher Ergebnisse ergeben.

Formal lässt sich Differential Privacy auf folgende Weise definieren:

### Definition 14.2.1

#### ( $\epsilon$ -Differential Privacy [28, Definition 1])

Eine randomisierte Funktion  $K$  stellt  $\epsilon$ -Differential Privacy sicher, wenn für alle Datensätze  $D_1$  und  $D_2$ , die sich um höchstens in einem Element unterscheiden, und alle  $S \subseteq \text{Range}(K)$  Folgendes gilt:

$$\mathbb{P}[K(D_1) \in S] \leq \exp(\epsilon) \times \mathbb{P}[K(D_2) \in S].$$

Die Wahrscheinlichkeit wird über die Münzwürfe von  $K$  bestimmt.

In diesem Zusammenhang würde die Ausgabe des Algorithmus  $K$ , der  $\epsilon$ -Differential Privacy bietet, bei einer Entfernung einer Zeile (Person) aus dem Datensatz ungefähr so wahrscheinlich bleiben, als ob die Zeile nicht entfernt würde. Daher bietet  $K$  Privatsphäre für jede einzelne Zeile, auch für den Fall, dass alle anderen Zeilen einem Angreifer bereits bekannt sind. Diese Definition wird immer wichtiger, da sie auf statistischen Grundlagen und Rechenleistung basiert – Art der Abfragen oder zusätzliche Nebeninformationen haben keine (signifikante) Wirkung, d. h. sie können das Niveau der Privatsphäre nicht reduzieren.

Der Parameter  $\epsilon$  kontrolliert dabei, wie groß der maximale Effekt eines Individuums auf das Ergebnis einer Analyse ist – man spricht deshalb von  $\epsilon$ -Differential Privacy. Im Umkehrschluss quantifiziert der Term  $\epsilon$  aber auch, inwieweit die Privatheit eines Individuums durch die Analyse kompromittiert werden kann. Kleinere  $\epsilon$ -Werte gehen daher mit einem höheren Grad an Privatheit, aber auch einer stärkeren Störung einher, was sich wiederum negativ auf die Qualität der Daten auswirkt [30].

Hier sieht sich der Anwender mit einem Interessenkonflikt konfrontiert – das Datensubjekt ist an einem möglichst umfassenden Schutz seiner persönlichen Informationen interessiert, der Anwender hingegen an aussagekräftigen, möglichst unverfälschten Daten. Dieser Konflikt wird zusätzlich dadurch erschwert, dass keine feste Richtlinie existiert, wie der Parameter  $\epsilon$  optimal gewählt werden soll. Keine analytischen Formeln oder statistische Vorarbeiten helfen bei der Wertbestimmung. Daher muss er manuell entsprechend der tatsächlichen Anwendung und der Sensitivität der Daten ausgewählt werden.

Die Wahl reicht von  $\epsilon = 0.01$  bis  $\epsilon = 1$  in akademischer Forschung bis hin zu Werten zwischen 1 und 10 in industriellen Anwendungsfällen (Google, Apple, US Census Bureau)[61].

Schließlich wird die Privatsphäre auch dann gewahrt, wenn eine Person mehr als eine einzige Zeile zur Datenbank beiträgt oder wenn eine Gruppe sich um ihre gemeinsamen Daten sorgt. Letzteres erfordert den Tausch von  $\exp(\epsilon)$  gegen  $\exp(\epsilon c)$ , wobei  $c$  die Größe der Gruppe angibt.

### 14.2.1 Differential Privacy für Frage-Antwort-Systeme

Differential Privacy wurde ursprünglich für Frage-Antwort-Systeme eingeführt, bei denen nur die aus den Originaldaten gewonnenen Antworten auf statistische Fragen an die Datenbank in anonymisierter Form herausgegeben werden sollen. Hier ist der Laplace-Mechanismus [29] der bekannteste Algorithmus. Dabei wird  $\epsilon$ -Differential Privacy durch das Hinzufügen einer bestimmten Menge an Rauschen zum eigentlichen Ergebnis der statistischen Abfrage erreicht.

Dies geschieht wie folgt: Sei  $f : D \rightarrow \mathbb{R}^k$  eine Funktion, die die Anfrage repräsentiert, wobei  $D$  hier die Menge aller zulässigen Datenbanken für diese Anfrage sei. Die wahre Antwort ist dann  $f(D)$  für eine Datenbank  $D \in D$ . Der Algorithmus  $K_f$  antwortet auf die Abfrage  $f$  als  $K_f(D) = f(D) + X$ , wobei  $X$  eine Zufallsvariable ist, die gemäß der Laplace-Verteilung  $\text{Lap}(\Delta f/\epsilon)$  verteilt ist, und  $\Delta f$  bezeichnet die Sensitivität von  $f$ , d. h. die maximale Differenz in den Werten der Antworten von  $f$  auf Paaren von Datenbanken, die sich in einem Element unterscheiden. Geht es beispielsweise um die Zählung von Personen mit bestimmten Merkmalen, so ist  $\Delta f = 1$ . Der auf diese Weise definierte Laplace-Mechanismus erfüllt Differential Privacy:

#### Theorem 14.2.1 (Dwork et al. [28, Theorem 1])

Für  $f : D \rightarrow \mathbb{R}^k$  erfüllt der Mechanismus  $K_f$ , welcher zu jedem der  $k$  Ausgabeterme unabhängig erzeugtes Rauschen mit Verteilung  $\text{Lap}(\Delta f/\epsilon)$  hinzufügt,  $\epsilon$ -Differential Privacy.

Da in der Regel eine Datenbank für mehr als eine Anfrage genutzt wird, gibt es in der Praxis ein Privacy-Budget, welches das über alle Antworten kumulierte  $\epsilon$  ist. Im einfachsten Fall gibt es ein fixes  $\epsilon$  für alle Antworten, sodass das Privacy-Budget nach einer festen Anzahl aufgebraucht ist. Wird eine Datenbank für eine festgelegte Analyse verwendet, so lässt sich die Menge der Anfragen vorher abschätzen, sodass das Budget passend verteilt werden kann. Ist die Anzahl der Anfragen vorab unbekannt, etwa bei einem öffentlichen Auskunftssystem, wird bei festem  $\epsilon$  das Privacy-Budget nach und nach von den Antworten aufgebraucht, sodass irgendwann gar keine Antworten mehr gegeben werden können, wenn das Budget erschöpft ist. Ein Ausweg ist hier eine sukzessive Reduktion von  $\epsilon$  nach einem geeigneten mathematischen Schema, sodass das Budget nie vollständig aufgebraucht wird. Dies bedeutet aber auch, dass die Störung wird sukzessive erhöht wird und der Informationsgehalt von Antworten dadurch sukzessive reduziert wird, sodass die Antworten irgendwann unbrauchbar werden.

---

### 14.2.2 Weitere Anwendungen für Differential Privacy

---

Neben dem Frage-Antwort-System gibt es noch weitere Anwendungsszenarien für Differential Privacy. Der Exponential-Mechanismus [55] kann genutzt werden, um synthetische Tabellen nach dem Vorbild der Originaltabelle zu erzeugen [9]. Solche synthetischen Daten sind beispielsweise zum Testen von Software geeignet, aber auch für Anwendungen des maschinellen Lernens, die auf privatsphärenfreundliche Weise statistische Eigenschaften der Originaldaten erlernen sollen (vgl. Kapitel 16). Der Exponential-Mechanismus begünstigt zum einen die Nähe zu den Originaldaten mit höheren Wahrscheinlichkeiten und erfüllt zum anderen das Kriterium von Differential Privacy durch eine Zufallsstreuung, deren Größe über einen Parameter gesteuert wird. Der Exponential-Mechanismus ist jedoch mit sehr hohem Aufwand verbunden.

Ein weiteres Verfahren für Differential Privacy sind randomisierte Antworten, welche bereits lange in sozialwissenschaftlichen Studien eingesetzt werden [85]. Dabei geben Probanden in Abhängigkeit von Münzwürfen oder Ähnlichem zufallsbestimmte oder wahrheitsgemäße Antworten. So lässt sich aus der einzelnen Antwort keine Wahrheit ablesen, d.h.

die Privatsphäre der Probanden wird schon bei der Datenerhebung geschützt. Durch das Gesetz der großen Zahlen kann der Einfluss der Zufallsantworten auf die Gesamtheit der Antworten näherungsweise herausgerechnet werden, sodass mit statistischen Methoden Erkenntnisse aus den Daten abgeleitet werden können.

Nachdem das Konzept von Differential Privacy erfunden war, wurde nachgewiesen, dass randomisierte Antworten in der Tat die Anforderungen von Differential Privacy bei geeignetem Studienaufbau erfüllen [41]. Dadurch ist ein effektiver Schutz der Privatsphäre gegeben. Zudem ist das Verfahren rechnerisch (aber nicht unbedingt für die Probanden) effizient. Randomisierte Antworten stellen sogar Local Privacy sicher, so dass kein Vertrauen in den Datenaggregator nötig ist [25]. Zu beachten bleibt aber, dass bei randomisierten Antworten ein deutlicher Informationsverlust entsteht, und dass die Daten in ihren statistischen Eigenschaften hochgradig verändert werden. Letzteres kann rechnerisch korrigiert werden, aber ersteres kann nur mit einer größeren Probandenzahl kompensiert werden.

## 15 Anonymisierung von Texten

Bei Textdokumenten unterscheiden wir zwischen der Meta-datenebene, der Inhaltsebene und der Schreibstilebene. Auf all diesen Ebenen können Personenbezüge vorhanden sein. Bevor wir auf die Anonymisierung hinsichtlich jeder einzelnen Ebene eingehen, klären wir zunächst diese Begriffe. Die Meta-datenebene ist eine vom Text entkoppelte Ebene, die Zusatzinformationen zu einem Dokument bereitstellt. Die Inhaltsebene ist die zentrale Ebene, die die eigentliche Information trägt. Die Schreibstilebene ist in die Inhaltsebene eingebettet und

lässt sich nicht ohne Weiteres von dieser entkoppeln.

Verräterische Spuren hinsichtlich der Identität der Person(en), die das Dokument erstellt haben, oder Informationen bzgl. Dritter können sich in diesen Ebenen wiederfinden. Im Folgenden gehen wir darauf ein, wie eine Anonymisierung auf allen Ebenen ermöglicht werden kann, unter der Anforderung, dass der Nutzen der Daten möglichst erhalten bleibt.

### 15.1 Anonymisierung auf der Metadatenebene

Die Existenz und Form von Metadaten hängt davon ab, in welchem Format ein Dokument vorliegt. Handelt es sich um eine Datei in einem komplexen Format (z. B. eine PDF-Datei oder ein Word-Dokument), so liegen in der Regel Metadaten vor. Diese enthalten Felder wie etwa Autoren, Titel, Schlüsselwörter und Erstellungsdatum und reichern das Dokument mit semantischen Informationen an. Handelt es sich jedoch um eine reine Textdatei, so existiert innerhalb der Datei keine Metadatenebene. Gegebenenfalls finden sich jedoch Metadaten im umgebenden System, welches die Datei speichert, was beispielsweise ein Dateisystem oder eine E-Mail sein kann.

Metadaten bergen die Gefahr, dass sie oft vom Ersteller nicht wahrgenommen werden, jedoch Informationen enthalten, die dessen Identität ungewollt preisgeben können. Dies kann

je nach Kontext unterschiedliche Konsequenzen haben. Ein mögliches Szenario wäre z. B. die Einreichung einer Publikation bei einer Konferenz, die Doppelblindgutachten (engl. double-blind review) durchführt. Hierbei muss seitens des einreichenden Autors gewährleistet sein, dass keinerlei Informationen bzgl. dessen Identität preisgegeben werden, damit die Gutachter ein unvoreingenommenes Gutachten durchführen können.

Die Anonymisierung der Metadatenebene ist meist trivial durchführbar, indem die Metadaten entweder gar nicht erst erstellt oder nachträglich entfernt werden. In beiden Fällen bleibt der Inhalt des Dokuments unversehrt, da die Metadaten, wie oben erwähnt, von dem eigentlichen Text entkoppelt sind.

### 15.2 Anonymisierung auf der Inhaltsebene

Inhaltsdaten enthalten oftmals Entitäten wie z. B. Personennamen, Bezeichnungen von Firmen oder Organisationen oder geografische Orte, die die Identität des Autors oder die von Dritten referenzieren können. Diese lassen sich anders als Metadaten nicht mit einfachen Mitteln entfernen,<sup>1</sup> ohne die Semantik des Dokuments zu verletzen. Die Voraussetzung für die Anonymisierung von Texten ist, zunächst die Verweise auf Identitäten zu identifizieren. Diese können mithilfe computerlinguistischer Verfahren wie Eigennamenerkennung (engl. na-

med entity recognition) ermittelt werden [45, 86]. Anschließend können diese Verweise mit verschiedenen Strategien anonymisiert werden. Eine hundertprozentige Erkennung aller Verweise ist jedoch nicht möglich, sodass immer ein Restrisiko verbleibt.

Eine Möglichkeit zur Anonymisierung entsprechender Textstellen läuft über eine Pseudonymisierung mittels partieller Verschlüsselung. Partielle Verschlüsselung wird in erster Linie

<sup>1</sup> Ausgenommen sind isolierte Entitäten, die unabhängig vom Text sind (z. B. der Name nach einer Grußformel).

auf Audiodaten [73], Bildern [18] oder Videos [42] eingesetzt, lässt sich jedoch auch auf Texte übertragen. Dabei werden Verweise auf Identitäten mit einem geheimen Schlüssel  $k$  verschlüsselt, sodass aus dem Dokument  $D$  ein modifiziertes Dokument  $D'$  entsteht.  $D'$  kann somit nur von autorisierten Personen, die  $k$  besitzen, entschlüsselt und dadurch vollständig gelesen werden. Stellt man sicher, dass nach der Pseudonymisierung niemand mehr den Schlüssel  $k$  hat oder herleiten kann, ist eine Anonymisierung erreicht. Der Nachteil der partiellen Verschlüsselung ist, dass der Lesefluss in  $D'$  durch die verschlüsselten Elemente gestört wird und das Dokument daher nur fragmentarisch gelesen werden kann, was den Nutzen des Dokuments reduziert.

Eine Alternative zur partiellen Verschlüsselung ist, die Verweise auf Identitäten zu paraphrasieren. Damit kann eine Anonymisierung erreicht und gleichzeitig die Semantik von  $D$  bis zu einem gewissen Grad beibehalten werden. Analog zur partiellen Verschlüsselung führt dies zwar ebenfalls zu einem Informationsverlust, allerdings in einer Form, bei der zum einen die modifizierte Version  $D'$  vollständig lesbar bleibt und zum anderen niemand mit einer Art Schlüssel die Ursprungsinformation wiederherstellen kann. Dazu gilt es, die identifizierten Entitäten durch generischere Angaben<sup>2</sup> zu ersetzen. Anhand von Sprachmodellen (engl. language models) an, die sonst insbesondere im Bereich der maschinellen Übersetzung Anwendung finden, kann mitunter bewertet werden, ob modifizierte Phrasen natürlich wirken oder in der jeweiligen Sprache eher untypisch sind.

Eine wichtige Frage bei der Paraphrasierung ist, woher die abgewandelten Entitäten bezogen werden können. Eine Möglichkeit besteht darin, vorhandene linguistische Ressourcen zu verwenden, wie etwa Ontologien oder lexikalische Wortnetze, mit denen semantisch sinnvolle Ersetzungen durchgeführt werden können. Diese müssen in der Regel händisch erstellt werden und sind dadurch mit entsprechendem Aufwand und hohen Kosten verbunden. Hinzu kommt die Problematik der temporalen Veränderung von Sprachen,<sup>3</sup> sodass gegebenenfalls zu einer Entität  $x$  in einem Text keine passenden Ersetzungen in einer Wortliste gefunden werden können, da die Wortliste zu einem Zeitpunkt erstellt wurde, als  $x$  noch nicht existierte.

Alternativ zu händisch erstellten linguistischen Ressourcen eignen sich Ansätze basierend auf sogenannte Word Embeddings. Die Idee dahinter ist, Wörter eines Vokabulars als reelle Vektoren in einem hochdimensionalen Raum darzustellen und diesen auf einen Raum mit niedrigerer Dimension abzubilden, sodass im zweiten Raum semantische Beziehungen der Wörter durch die Nähe der entsprechenden Vektoren widergespiegelt werden. Verwendet wird dafür ein Maß (z. B. die Kosinus-Ähnlichkeit), das Wort-Vektoren als Eingabe erhält und einen Wert im Bereich  $[0; 1]$  liefert, welcher angibt, wie ähnlich beide Wörter (bzw. deren Repräsentationen) zueinander sind. Mithilfe solcher Word Embeddings lassen sich ohne den Einsatz gelabelter Daten bzgl. einer Entität  $x$  semantisch ähnliche Entitäten  $y_1, y_2, \dots$  finden, die eine Ersetzung erlauben. Vorausgesetzt werden hier jedoch genügend ungelabelte Textdaten, welche Informationen über die Entität  $x$  enthalten. Ein wesentlicher Nachteil hierbei ist allerdings, dass die Entitäten nicht in einer festgelegten Relation (z. B. Synonymie) zueinander stehen, sondern sich über mehrere Relationen wie etwa Hyperonymie, Hyponymie, Meronymie oder Holonymie erstrecken können. Der Literatur zufolge existiert noch kein zufriedenstellender Ansatz, mit dessen Hilfe Entitäten hinsichtlich ihrer semantischen Relationen automatisiert abgegrenzt werden können, sodass es hierfür noch weiterer Forschungsarbeit bedarf. Als mögliche Lösung dazu wurden in den letzten Jahren sogenannte Sense Embeddings [38] entwickelt, mit deren Hilfe sich ähnliche Wörter hinsichtlich ihrer semantischen Relation eingrenzen lassen, sodass für  $x$  z. B. nur noch Hyperonyme bestimmt werden, welche eine Anonymisierung gewährleisten. In dieser Hinsicht sind Sense Embeddings eine vielversprechende Option, um eine Anonymisierung auf der Inhaltsebene zu ermöglichen. Allerdings werden auch hier zum Lernen der relevanten Relationen nicht-anonyme Textdaten benötigt.

Eine weitere Strategie zur Paraphrasierung basiert auf Koreferenzen. In einem Text bezeichnet man mit Koreferenzen Verweise auf eine zuvor bereits aufgetretene Entität. Dies kann über Wiederholungen der Bezeichnung der Entität, über Pronomen oder über Umschreibungen geschehen. Letztere Art von Koreferenzen sind eine geeignete Quelle von Paraphrasen. Bei einer Anonymisierung über Koreferenzen gilt es also, nach der Erkennung von Entitäten die Koreferenzen zu erkennen und in einem weiteren Schritt identifizierte Bezeich-

2 Beispielsweise „Angela Merkel“  $\rightarrow$  { „deutsche Politikerin“, „gebürtige Hamburgerin“, ... }.

3 Vor 20 Jahren gab es z. B. noch nicht die Wörter „googlen“, „Podcast“ und „Smombie“.



nungen durch generischere Angaben zu ersetzen, die aus Koreferenzen gewonnen werden. Eine initiale Untersuchung und Erprobung dieser Strategie hat Roey Regev in seiner Ba-

chelorarbeit [65] an der TU Darmstadt in Kooperation mit dem Fraunhofer SIT durchgeführt.

## 15.3 Anonymisierung auf der Schreibstilebene

Die Identität einer Person lässt sich auch über dessen Schreibstil bestimmen. Im Laufe des letzten Jahrzehnts hat sich die digitale Textforensik als Forschungsfeld etabliert, welches sich mit der Authentizität und Glaubwürdigkeit von Texten auseinandersetzt. Hauptaugenmerk liegt dabei auf der Autorschaftsanalyse, welche das Ziel verfolgt, Informationen über die Autoren digitaler Dokumente offenzulegen [64]. Zu den bekanntesten Disziplinen dieses Forschungsfelds zählen die Autorschaftsattributions (AA) und die Autorschaftsverifikation (AV).

In der AA gilt es, hinsichtlich eines Dokuments  $D_U$  mit unbekannter Autorschaft und einer Menge von potenziellen Kandidaten  $\{A_1, A_2, \dots\}$ , von denen jeweils Beispieltexthe vorliegen, den wahrscheinlichsten Autor  $A_i$  von  $D_U$  zu identifizieren. In der AV gilt es hingegen, für  $D_U$  und eine Referenzmenge von Texten nur eines bekannten Autors  $A$  zu entscheiden, ob  $D_U$  von  $A$  verfasst wurde. In der Literatur wurden bezüglich beider Disziplinen zahlreiche Ansätze vorgeschlagen, die vielversprechende Ergebnisse, auch über mehrere Sprachen hinweg, hervorgebracht haben.

Ursprünglich wurden AA/AV-Verfahren mit der Vision entwickelt, sie in forensischen Szenarien einzusetzen (z. B. um zu beantworten, ob mehrere Erpresserschreiben von ein und demselben Täter geschrieben wurden). Allerdings wurde im Laufe der Zeit auch ersichtlich, dass sich mit den stetig wachsenden Erkennungsleistungen der AA/AV-Verfahren sich die eigentliche Anwendung zweckentfremden lässt, um die Anonymität beliebiger (nicht-krimineller) Personen auch jenseits von forensischen Szenarien zu demaskieren.

Aus der Notwendigkeit heraus, die Identität von Autoren zu schützen, entstand das Forschungsfeld Author Obfuscation (AO), welches sich damit befasst, wie sich der Schreibstil in Dokumenten verschleiern lässt und damit die Identität der Autoren anonymisieren lässt. Bisherige AO-Ansätze lassen sich in manuelle, computerassistierte und automatische Verfahren aufteilen [37], wobei der Forschungsfokus insbesondere auf letzteren liegt. Automatische AO gilt als sehr anspruchsvoll, da sie auf Sprachkompetenzen zurückgreifen muss, um anonymi-

sierende Umformungen in den Dokumenten vorzunehmen, bei gleichzeitiger Beibehaltung der ursprünglichen Semantik.

Unter den veröffentlichten automatischen AO-Verfahren ist vor allem der Ansatz Adversarial Author Attribute Anonymity Neural Translation (A<sup>4</sup>NT) von Shetty et al. [74] hervorzuheben. Das Verfahren ist unserer Recherche nach das einzige Verfahren, das eine dedizierte Komponente für die Semantikerhaltung enthält. A<sup>4</sup>NT verfolgt eine intuitive Idee, die analog zu einer maschinellen Übersetzung funktioniert. Während in der maschinellen Übersetzung ein Dokument in eine festgelegte Zielsprache übersetzt wird, wird bei A<sup>4</sup>NT das Dokument in dieselbe Sprache wie die Quellsprache „übersetzt“, jedoch in einer identitätsverschleiernenden Form. Dadurch soll z. B. eine Täuschung von AA/AV-Verfahren erreicht werden, um den Schreibstil des ursprünglichen Autors nicht mehr wiedererkennen zu können. Um diese Texttransformation durchzuführen, verwendet A<sup>4</sup>NT mehrere neuronale Netze, wobei der Kern ein sogenanntes GAN-Netzwerk ist. Dieses besteht wiederum aus einem Klassifikator  $K$  und einem Generator  $G$ . Die Aufgabe von  $K$  ist es, zwischen zwei oder mehr Klassen (z. B. verschiedene Autoren, oder autorspezifische Eigenschaften wie Geschlecht oder Alter) zu unterscheiden, während  $G$  versucht,  $K$  zu täuschen und dabei gleichzeitig die Semantik des Originaltextes beizubehalten. Die Komponente, die für die Semantikerhaltung zuständig ist, basiert auf einer Kombination zweier Maße: Die Wahrscheinlichkeit der Rekonstruktion des Originaltextes mithilfe einer Rücktransformation seitens A<sup>4</sup>NT und die Distanz bzgl. der Sentence Embeddings. Das Verfahren wurde hinsichtlich der drei autorspezifischen Attribute Alter (unter 20 vs. über 20), Geschlecht und Identität (Obama vs. Trump) anhand einer Kollektion von Blogartikeln und einer Kollektion von politischen Reden getestet. Shetty et al. konnten die Erkennungsgenauigkeit (F1-Wert) beim Alter von 88 % auf 8 %, beim Geschlecht von 75 % auf 39 % und bei der Identität von 100 % auf 0 % senken, was dafür spricht, dass eine Anonymisierung auf der Schreibstilebene möglich ist. Genauer gesagt ist eine (bei Alter und Identität) mehr oder (beim Geschlecht) weniger zuverlässige Imitation der jeweiligen Gegenklasse gelungen.

## 16 Anonymisierung im Kontext von maschinellem Lernen

Die Kombination aus verhältnismäßig günstigem, skalierbarem Speicherplatz, leistungsstarken Rechenkapazitäten und dem Bestreben, wertvolle Informationen aus den Daten, die kontinuierlich generiert werden, zu gewinnen, verhalf dem maschinellen Lernen zu seinem heutigen Erfolg. Maschinelles Lernen (ML) ist ein Teilgebiet der künstlichen Intelligenz und beschreibt eine Reihe von Lernalgorithmen, die versuchen, Strukturen in Daten zu erkennen, um basierend auf diesen Mustern bspw. Klassifizierungs- oder Regressionsaufgaben zu lösen. Der Einsatz von Verfahren des maschinellen Lernens bietet sich immer dann an, wenn die zu lösenden Probleme zu

komplex oder zu umfassend sind, um sie analytisch beschreiben zu können [21]. Gleichzeitig bedeuten größere Datenmengen auch, dass mehr Informationen zum Trainieren der Lernalgorithmen zur Verfügung stehen, was tendenziell zu besseren Modellen und effizienteren Schätzungen führt [77]. Neuronale Netze finden aufgrund ihrer Flexibilität und guten Generalisierungsfähigkeit in den verschiedensten Bereichen Anwendung – ob im Verarbeiten und Analysieren natürlicher Sprachen, zur Bild- oder Gesichtserkennung oder zum Aufspüren von Anomalien.

### 16.1 Privatsphärenrisiken beim maschinellen Lernen

Im Zeitalter von Big Data und maschinellem Lernen (ML) ist es noch schwieriger geworden, Privatheit zu gewährleisten, da in großen Datenbeständen – selbst in solchen aus gering strukturierten oder gar unstrukturierten Daten – die entschei-

denden Verknüpfungen gefunden werden können, welche das Herstellen von Personenbezügen ermöglichen. Wie wenig Information benötigt wird, um ein Individuum eindeutig identifizieren zu können, zeigen die Beispiele aus Kapitel 3.

---

#### 16.1.1 Risiken durch gelernte Modelle

---

Da ML-Algorithmen üblicherweise auf disjunkten Datensätzen trainiert und evaluiert werden, wurde lange fälschlicherweise angenommen, dass es nicht möglich ist, vom finalen Modell Rückschlüsse auf die zum Training verwendeten Daten zu ziehen, was folglich einer Anonymisierung des verwendeten Datenmaterials gleichkommen würde. Bestimmte ML-Techniken können sich jedoch unerwartet deutlich an die zum Training des Modells verwendeten Daten erinnern. So speichern Support Vector Machines oder k-nächste-Nachbarn-Klassifikatoren Informationen über die zum Lernen verwendeten Daten in dem Modell selbst ab. Diese sogenannten Feature-Vektoren erlauben unter bestimmten Umständen Rückschlüsse auf die Rohdaten und stellen somit ein entscheidendes Risiko dar [68].

Fredrikson et al. [33] demonstrierten, dass die Erinnerung in neuronalen Netzen, welche zur Gesichtserkennung genutzt wurden, mitunter so stark sein kann, dass es möglich ist, ein Abbild der Trainingsdaten zu rekonstruieren – ein sogenannter

Modellinversionsangriff. Shokri et al. [76] bewiesen, dass neuronale Netze aufgrund ihrer Konstruktion anfällig für Membership-Inference-Angriffe sind. Die Autoren wiesen nach, dass ein trainiertes Netz merkbar anders auf Informationen reagiert, welche bereits zum Training verwendet wurden, als auf bisher ungesehene Testdaten. Aufgrund dieser Rückmeldung kann ein Angreifer eindeutig zuordnen, ob ein Individuum in einem bestimmten Datensatz enthalten ist oder nicht. Solche Angriffe stellen allgemein eine Verletzung der Privatheit dar, sind aber besonders dann kritisch, wenn es sich um sensible Informationen handelt, wie beispielsweise Insolvenz oder ob eine bestimmte Krankheit vorliegt.

---

### 16.1.2 Risiken durch maschinelles Lernen in der Cloud

---

Die allermeisten potenziellen Anwender von maschinellem Lernen sehen sich mit mindestens einem der folgenden drei Probleme konfrontiert: Ihnen fehlt es an nötigem Fachwissen, sie verfügen nicht über hinreichende Rechenkapazitäten oder ihre Datenbasis ist nicht ausreichend, um ihre Modelle gut trainieren zu können. Sie sind folglich darauf angewiesen, ihre Daten in irgendeiner Weise zu teilen. Daher werden fortschrittliche Analysen häufig in Kombination mit Cloud Computing angeboten, welches jedoch grundsätzlich ein Sicherheitsrisiko birgt, da der Datenbesitzer in aller Regel nicht deckungsgleich mit der Partei ist, die den Service – z.B. Machine Learning as a Service (MLaaS) – bereitstellt. Oft werden die Daten, die an den Server übermittelt wurden, unverschlüsselt auf diesem gelagert, was sie nicht nur anfällig für Missbrauch von Seiten des Serviceproviders macht, sondern ebenfalls für Angriffe unbefugter Dritter auf die Server.

Vor dem Hintergrund der Vielzahl an Datenskandalen in den vergangenen Jahren im Allgemeinen und der oben aufgezeigten Angriffe auf maschinelles Lernen im Speziellen mag es ratsam erscheinen, Daten aus dem eigenen Kontrollbereich gar nicht erst herauszugeben, um Privatheit zu schützen. Dies hätte jedoch einschneidende Konsequenzen für die Forschung, beispielsweise für die Medizinforschung. Außerdem könnten viele weitverbreitete, nützliche Dienste nicht weiter angeboten werden. Auch ML-Systeme zur Strafverfolgung, etwa zur automatisierten Aufdeckung von Kinderpornografie, können ohne die Bereitstellung von geeigneten Trainingsdaten nicht trainiert und somit nicht eingesetzt werden. Ziel ist es folglich, Daten auf privatsphärenfreundliche Weise einem ML-System zur Verfügung stellen zu können.

## 16.2 Privatsphärenfreundliches maschinelles Lernen

Das Forschungsfeld Privacy-preserving Machine Learning (PPML) ist noch recht jung. Das Ziel ist, die Privatheit des Einzelnen zu schützen und gleichzeitig das Training von Modellen auf Daten von vielen Personen zu ermöglichen. Auch wenn es bereits vielversprechende Ansätze gibt, besteht noch viel Entwicklungsbedarf. Nachfolgend werden die wichtigsten Forschungsrichtungen auf diesem Gebiet skizziert. Möglichkeiten zum Schutz sensibler Daten vor der Offenlegung beim maschinellen Lernen bieten u. a. das Konzept von Differential Privacy, Verfahren der homomorphen Verschlüsselung und das kollaborative maschinelle Lernen.

In Bezug auf die Erkennung von Kinderpornografie haben die Zentral- und Ansprechstelle Cybercrime (ZAC) aus Nordrhein-Westfalen und Microsoft Deutschland im August 2019 eine Kooperation bekanntgegeben, in der ein Ansatz zum privatsphärenfreundlichen Training eines ML-Systems erprobt werden soll.<sup>1</sup> Dabei verlassen die Rohdaten nicht die Strafverfolgungsbehörde. Stattdessen extrahiert die Behörde bestimmte Feature-Vektoren aus den einschlägigen Bildern und geben diese abstrakten Daten in die Microsoft-Cloud, wo sie zum Training des ML-Systems verwendet werden oder bewertet werden. Dies stellt eine Vorstufe des kollaborativen maschinellen Lernens dar.

---

<sup>1</sup> <https://news.microsoft.com/de-de/ki-im-einsatz-gegen-kinderpornografie/>.

### 16.2.1 Differential Privacy für maschinelles Lernen

Arbeiten zu Differential Privacy (vgl. Abschnitt 14.2) im Kontext des maschinellen Lernens erforschen verschiedene Aspekte des Verrauschens von potenziell angreifbaren Daten. Untersucht wird hier meist, auf welcher Ebene die Störungen idealerweise Eingang in den Algorithmus finden – ob nun auf Input- oder Output-Ebene, oder ob die Gradienten oder die Verlustfunktion verrauscht werden – und welche Verteilungseigenschaften das Rauschen selbst haben sollte. Das Ziel ist, einen optimalen Trade-off zwischen Privatheit und Ergebnisqualität zu erreichen.

Eine andere Richtung der privatheiterhaltenden Datenveröffentlichung und -analyse verfolgt der Ansatz der Differentially

Private Data Synthesis (DIPS). Hierbei werden Daten auf Basis realer Datensätze beispielsweise mittels Copula-Funktionen [44] oder Generative Adversarial Networks [83] unter Einhaltung von Differential Privacy synthetisiert. Der offensichtliche Vorteil dieses Ansatzes ist, dass die simulierten Daten bereits Differential Privacy erfüllen und somit keine Rückschlüsse auf die Ursprungsdaten ermöglichen – im Gegensatz zu anderen Datensyntheseverfahren. Darüber hinaus besitzen die Daten annähernd die gleichen Verteilungseigenschaften wie die zugrunde liegenden Originaldaten und können in beliebiger Anzahl generiert werden, um so beispielsweise die Güte eines ML-Modells zu verbessern [53, 61].

### 16.2.2 Maschinelles Lernen mit homomorpher Verschlüsselung

Homomorphe Verschlüsselung (vgl. Abschnitt 13.1) erlaubt – im Gegensatz zu herkömmlichen Verschlüsselungsmethoden – Rechenoperationen direkt auf den verschlüsselten Daten auszuführen, ohne diese zuvor in Klartext überführen zu müssen und sie dadurch angreifbar zu machen. Jede Operation liefert ein ebenfalls verschlüsseltes Ergebnis, das dechiffriert demjenigen entspricht, welches resultieren würde, wäre die Operation auf dem entsprechenden Klartext durchgeführt worden. Mit homomorpher Verschlüsselung können daher Daten an eine nicht-vertrauenswürdige Instanz weitergegeben werden und Berechnungen dort durchgeführt werden.

Insbesondere die sogenannte voll-homomorphe Verschlüsselung generiert jedoch einen signifikanten Rechenmehraufwand [23, 49], der diese für rechenintensive Anwendungen wie maschinelles Lernen bisher unbrauchbar macht. Erste praktikable Ansätze verwenden daher Vereinfachungen. So wenden Dowlin et al. [23] ein auf unverschlüsselten Rohdaten trainiertes neuronales Netz auf begrenzt (engl. somewhat) homomorph verschlüsselte Daten an. Long et al. [50] hingegen haben das Training verschiedener ML-Verfahren mit additiv-homomorpher Verschlüsselung und Zero-Knowledge-Beweisen realisiert.

### 16.2.3 Kollaboratives maschinelles Lernen

Eine Lösung zum privatsphärenfreundlichen Lernen auf Daten von vielen Nutzern ist das kollaborative maschinelle Lernen, auch bezeichnet als verteiltes, dezentrales oder föderiertes Lernen. Hierbei trainieren die Nutzer ein Grundmodell lokal auf ihren individuellen Daten und übermitteln lediglich die neu berechneten Gradienten des Trainings oder die neuen Modellparameter an den Serviceprovider. In einem periodischen Prozess aktualisiert der Provider das Gesamtmodell anhand der übermittelten Informationen aller Teilnehmer und stellt es ihnen anschließend zum Download zur Verfügung. Diese trainieren nun das aktualisierte Modell erneut lokal und senden die resultierenden Gradienten oder Parameter zurück an den Server [54].

Hauptsächlich zum Schutz der Privatsphäre, aber auch zur Kommunikationseffizienz, erlaubt der Ansatz nach Shokri und Shmatikov [75], dass nicht alle Aktualisierungen mit dem Server geteilt werden müssen, sondern nur eine kleine Teilmenge, deren Größe vom Nutzer selbst festgelegt wird. Zusätzlich wird in diesem Ansatz das Konzept von Differential Privacy (vgl. Abschnitt 14.2 sowie unten) umgesetzt. Allerdings sollte sich der Nutzer des Trade-offs zwischen der Menge der geteilten Aktualisierungen sowie Trainingszeit und -qualität bewusst sein.

Bei der Umsetzung von kollaborativem Lernen können auch kryptografische Methoden genutzt werden, etwa homomorphe Verschlüsselung (s. o.) oder die sichere Mehrparteienberechnung (vgl. Abschnitt 13.2). Das Ziel von MPC ist das gemeinschaftliche Berechnen einer Funktion, für die mehrere Parteien eine Eingabe liefern. Die Privatheit wird in dieser Art der Berechnung dadurch gewahrt, dass jede der beteiligten Parteien nur das Endergebnis, d. h. die Funktionsausgabe, und die eigene Eingabe erfährt. Die Eingaben der übrigen Teilnehmer bleiben verborgen. Je nach Anzahl der Teilnehmer und deren Abbruchwahrscheinlichkeit existieren verschiedene Ansätze mit unterschiedlichem Rechen- und Kommunikationsaufwand, um dieses Ziel zu erreichen. So kann die Berechnung der Funktion bspw. auf mehrere nicht-kolludierende Server aufgeteilt werden. Tatsächlich gibt es erste MPC-Ansätze im Kontext von maschinellem Lernen zur Summenberechnung von Modellparametern [10].

Hitaj et al. [39] haben dargelegt, dass es selbst in dezentralen Lernansätzen mit Hilfe eines Generative Adversarial Networks (GAN) möglich ist, über die übrigen aufrichtigen Teilnehmer sensible Daten zu sammeln. Der Angriff von Hitaj et al. ist sogar durchführbar, wenn das dezentrale Lernen mit Differential Privacy wie bei Shokri und Shmatikov [75] oder mit MPC wie bei Bonawitz et al. [10] kombiniert wird. Melis et al. [56] entkräften teilweise die Argumente von Hitaj et al., zeigen aber selbst neue Angriffsstrategien auf.

# 17 Zu bewältigende Herausforderungen

In diesem Kapitel möchten wir einen Blick in die Zukunft werfen: Was muss getan werden, um die aktuelle Technik so zu verbessern, dass Privatsphärenschutz im Kontext von Big Data besser gewährleistet und angewendet werden kann? Wir stellen einen Überblick von zu verbessernden Konzepten und Methoden sowie von zu schaffenden nicht-technischen Rahmenbedingungen zur Verfügung, die die Implementierung einer datenschutzgerechten bzw. privatsphärenhaltenden Verarbeitung von Big Data unterstützen sollen. Die ersten Abschnitte konzentrieren sich auf die nahe Zukunft und decken

Ansätze ab, die heute mit Einschränkungen genutzt werden können und die in den nächsten fünf Jahren voraussichtlich weiterentwickelt werden. Wir weisen auf spezifische Aspekte hin, die es zu verbessern gilt. Abschnitt 17.4 konzentriert sich auf Aspekte, die eine Entwicklung von mindestens einem Jahrzehnt erfordern, bis sich Auswirkungen auf die Praxis ergeben. Wir diskutieren insbesondere verschiedene Aspekte zum Thema Anonymisierung mit kurzem und langem zeitlichen Horizont.

## 17.1 Technische Fortschritte

Die folgenden Unterabschnitte zeigen, wo technologische Verbesserungen in Bezug auf den Datenschutz dringend nötig sind.



Abbildung 17.1: Die Zukunft von Big Data muss die Privatheit erreichen.

### 17.1.1 Effektive und effiziente Anonymisierung

Die Unterscheidung zwischen personenbezogenen und nicht personenbezogenen Daten ist in der Praxis schwierig. Angeblich anonyme Daten haben sich in der Vergangenheit oft als personenbezogene Daten herausgestellt, vgl. Kapitel 3.

Angesichts der Leistungsfähigkeit von Big-Data-Systemen ist es noch schwieriger geworden, Privatheit zu gewährleisten, da in großen Datenbeständen – selbst in solchen aus gering strukturierten oder gar unstrukturierten Daten – die entscheidenden Verknüpfungen gefunden werden können, welche das Herstellen von Personenbezügen ermöglichen. Daher genügen einfache, heuristische Anonymisierungen, wie sie im letzten Jahrhundert Stand der Technik waren, bei weitem nicht, sondern es sind Lösungen zur nachweisbaren Anonymisierung nötig, die selbst vom heutigen Stand der Technik nicht geliefert werden.

Bestehende Anonymitätskriterien (vgl. Kapitel 14) schützen nur gegen bestimmte Risiken und auch nur, wenn die Annahmen über das Hintergrundwissen der Angreifer und über die Eigenschaften der Daten korrekt sind. Daher scheitern sie oft daran, dass die zugrunde liegenden Annahmen im Angreifermodell nicht der Realität entsprechen oder Rückschlussmöglichkeiten innerhalb dieser Modelle übersehen werden. Dies gilt insbesondere für Big-Data-Systeme mit Zugriff auf umfangreiche Datenquellen.

Um das Problem von nicht berücksichtigten Angriffsmöglichkeiten zu lösen, muss die Forschung entweder ultimative Anonymitätskriterien finden oder sie muss wenigstens Anwender dabei unterstützen, Schutzlücken oder unpassende Annahmen aufzudecken. Ersteres lässt sich wohl nur langfristig lösen (vgl. Unterabschnitt 17.4.2). Für zweiteres sollten die existierenden Anonymitätskriterien durch formale Angreifermodelle ergänzt werden, welche die angenommenen Fähigkeiten und das angenommene (Hintergrund-)Wissen von Angreifern explizit machen. Anwender können damit leichter erkennen, was in ihren Szenarien durch welche Anonymisierung tatsächlich geschützt wird. Zusätzlich muss der Anwender aber auch stets die Semantik der vorhandenen Datenattribute bei der Wahl der Anonymisierung beachten. Aber selbst wenn all dies korrekt berücksichtigt wird, stellt sich die Frage, welcher Anonymitätsgrad nach dem jeweiligen Maß ausreichend ist, um von Anonymität im Sinne eines nicht mehr vorhandenen Personenbezugs sprechen zu können. Dieser Aspekt erfordert

eine Präzisierung der rechtlichen Anforderungen, vgl. Unterabschnitt 17.2.2.

Zwei weitere Probleme ergeben sich aus den bestehenden Anonymisierungsmethoden im Zusammenhang mit Big Data:

1. Diese Methoden sind in vielen Fällen nicht mehr effizient genug, um einsetzbar zu sein.
2. Diese Methoden führen zu einem großen Informationsverlust bei hochdimensionalen Daten, sodass der Nutzen von korrekt anonymisierten hochdimensionalen Daten verschwindet.

Dies sind zwei entgegengesetzte Anforderungen. Zum einen sind Methoden, die in Bezug auf den Informationsverlust optimale Lösungen für die etablierten Anonymitätskriterien liefern, sehr rechenintensiv. Zum anderen können effiziente Verfahren keine Optimalität der Ergebnisse garantieren. Daher sind gute Heuristiken nötig, die es erlauben, mit relativ geringem Aufwand möglichst nahe an optimale Lösungen zu gelangen. Ideal wäre eine adaptive Ausbalancierung der entgegengesetzten Anforderungen von Effizienz des Algorithmus und Minimalität des Informationsverlusts in Bezug auf das gewählte Anonymitätskriterium. Grundlegender ist die Frage nach Anonymitätskriterien, die möglichst wenig Informationsverlust fordern, um Anonymität zu erreichen, vgl. Unterabschnitt 17.4.2.

Eine weitere Forderung aus der Praxis ist die Anonymisierung für das Streaming von Daten, da bestehende Methoden in der Regel für statische Datensätze konzipiert sind und daher nur in der sogenannten Batch-Schicht von Big-Data-Systemen eingesetzt werden können. Ebenfalls zu lösen ist die Anonymisierung von verteilten Daten, wie sie typischerweise in Big-Data-Systemen vorkommen.

Die Anonymisierung unstrukturierter Daten ist auch für viele Szenarien der Verarbeitung von Big Data wichtig. Die wichtigsten Arten von relevanten unstrukturierten Daten sind Text und Fotos oder Videos. Fotos und Videos werden in der Regel anonymisiert, indem Gesichter erkannt und verwischt werden. Dies ist jedoch möglicherweise nicht ausreichend, da die verschwommenen Gesichter noch genügend Informationen zur Erkennung einer Person enthalten können, und es ist möglich,

eine Person anhand anderer Merkmale wie Kleidung oder individueller Bewegungsmuster zu erkennen.

Anonymisierungslösungen für die Inhaltsebene unstrukturierter Textdaten basieren auf der Erkennung und Ersetzung von Named Entities (s. Abschnitt 15.2). Es gibt jedoch keine automatisch überprüfbaren Kriterien für die Anonymität von unstrukturiertem Text, sodass die Anonymität noch weniger sicher ist als bei strukturierten Daten. Daher ist eine Anonymitätsgarantie bei Texten noch weniger möglich als bei strukturierten Daten. Zudem gibt es nach wie vor die Herausforderung der zuverlässigen Erkennung von Entitäten sowie Herausforderungen in Bezug auf Umsetzbarkeit und Anwendbarkeit von Strategien zur Ersetzung dieser Entitäten. Dies beginnt mit der Schwierigkeit, Entitäten sicher erkennen zu können und wird fortgesetzt von der Tatsache, dass alle Methoden zur Ersetzung von Entitäten Nachteile haben, etwa den Bedarf an linguistischen Ressourcen zu den Entitäten oder an nicht-anonymen Trainingsdaten. In Bezug auf die Anonymisierung der Stilebene (vgl. Abschnitt 15.3) gibt es erste empirische Ergebnisse, die erfolversprechend sind, aber eine allgemeine Zuverlässigkeit kann noch nicht daraus geschlossen werden.

Der Privatsphärenschutz in Verbindung mit maschinellem Lernen ist ein noch recht junges und unerforschtes Thema, das erst vor wenigen Jahren verschiedene Risiken aufgedeckt hat. Erste Lösungsansätze, etwa in Verbindung mit Differential Privacy oder Kryptografie, beschränken sich hauptsächlich auf den Schutz additiver Operationen. Die meisten privatheiterhaltenden Verfahren in Kontext von ML sind nur für die Anwendung auf einen bestimmten Lernalgorithmus optimiert und auf andere ML-Verfahren schwer bis gar nicht übertragbar. Zudem stellt mangelnde Skalierbarkeit ein Hindernis für die Anwendung privatheiterhaltender Maßnahmen in der Praxis dar. Das Schützen sensibler Informationen generiert immer zusätzliche Kosten – entweder aufgrund von höherem Berechnungsaufwand, extrem langen Trainingszeiten oder weil der Nutzen der Daten bspw. durch zugefügtes Rauschen vermindert wird. In manchen Fällen fallen diese Kosten bisher sogar so groß aus, dass eine Anwendung in der Praxis nicht tragbar ist [68].

---

### 17.1.2 Berechnungen auf verschlüsselten Daten

---

Die Verarbeitung großer Datenmengen erfordert oft den Einsatz von Cloud Computing. Bei der Verarbeitung von personenbezogenen Daten ist ein hohes Maß an Sicherheit erforderlich. Die Weitergabe dieser Daten an einen Dritten, der den Cloud-Service bereitstellt, stellt ein potenzielles Risiko dar. Dieses Risiko kann dadurch umgangen werden, dass die Daten ausschließlich in verschlüsseltem Zustand weiterverarbeitet werden (vgl. Kapitel 13). Dies ist prinzipiell mit homomorpher Verschlüsselung (s. Abschnitt 13.1) möglich, aber ein solches Vorgehen führt nach heutigem Stand noch zu einem nicht unerheblichen Mehrbedarf an Speicher und Rechenleistung. Deshalb ist es zwingend notwendig, die verschiedenen Formen der homomorphen Verschlüsselung weiterzuentwickeln. Sicherheit während des gesamten Verarbeitungszyklus muss gewährleistet sein, ohne dem Benutzer gleichzeitig unzumutbare Kosten aufzubürden. Nicht das Erreichen dieses Ziels selbst, sondern eher das Finden eines Kompromisses zwischen Sicherheit und Kosten scheint jedoch realistisch.

Einsatz von Cloud Computing. Jedoch birgt die Weitergabe personenbezogener Daten an die Cloud-Service-Provider ein erhöhtes Risiko für Betroffene und Verantwortliche. Daher sind starke Sicherheitsgarantien erforderlich. Eine attraktive Lösung ist die Verarbeitung der Daten in einem verschlüsselten Zustand. Das langfristige Ziel für die Berechnung von verschlüsselten Daten sind effiziente voll-homomorphe Verschlüsselungsverfahren (s. Unterabschnitt 17.4.1). Nach heutigem Stand ist voll-homomorphe Verschlüsselung nicht geeignet, da sie zu einem nicht praxistauglichen Mehrbedarf an Speicher und Rechenleistung führt.

Kurz- und mittelfristige Lösungen für die Berechnung von verschlüsselten Daten müssen in andere Richtungen gehen. Bestehende Ansätze für eine effiziente Berechnung verschlüsselter Daten verwenden z. B. deterministische Verschlüsselungsverfahren. Solche Schemata erlauben etwa die Suche in der verschlüsselten Domäne. Es gibt auch ordnungserhaltende Verschlüsselungssysteme, die typische Datenbankoperationen wie das Sortieren von Daten und Bereichsauswahloperationen

Die Verarbeitung großer Datenmengen erfordert oft den



ermöglichen. Partiiell homomorphe Verschlüsselung erlaubt schließlich die Addition oder Multiplikation von verschlüsselten Zahlen. Ein praktisches, wenn auch schon wieder veraltetes Werkzeug für die Arbeit mit verschlüsselten Datenbanken ist CryptDB<sup>1</sup>. Mit dem Encrypted Big-Query Client<sup>2</sup> gab es sogar auch ein erstes solches Tool für eine Big-Data-Umgebung, konkret den BigQuery-Service, aber die Entwicklung an dem Encrypted BigQuery Client wurde wieder eingestellt. Auch Crypsis und Cuttlefish (vgl. Abschnitt 13.1) sind zwei Prototypen für dieses Anwendungsfeld. Solche Werkzeuge sollten generell für verteilte Datenbanken von Big-Data-Systemen verfügbar sein.

Der größte Nachteil der deterministischen Verschlüsselung ist ihre schwache Sicherheit, insbesondere wenn ordnungserhaltende Systeme verwendet werden. Wirklich sicherere Systeme dürfen nicht deterministisch sein. Der Bereich der partiiell

homomorphen und begrenzt homomorphen Verschlüsselung ist eine vielversprechende Richtung, um einige Berechnungen zu ermöglichen, während höhere Sicherheitsgarantien gewährleistet sind und ohne die Effizienzprobleme der voll-homorphen Verschlüsselung in Kauf nehmen zu müssen. Die existierenden Systeme für verschlüsselte Datenbanken greifen jedoch auf deterministische und ordnungserhaltende Verschlüsselung zurück, da in homomorphen Schemata jede Vergleichsoperation ein verschlüsseltes Ergebnis liefert und komplexe Operationen wie das Sortieren von Tabellen voll-homorphen Verschlüsselung erfordern würde. Zusätzlich zu den Fragen der Anwendbarkeit und Sicherheit der verschiedenen kryptografischen Strategien bestehen bei Systeme, die verschiedene Verschlüsselungsmethoden gleichzeitig einsetzen, zusätzliche Sicherheitsfragen, die sich aus Kombination der Methoden ergeben (s. Abschnitt 13.1).

---

### 17.1.3 Synthetische Daten

---

Eine Alternative zur Anonymisierung echter Daten stellt die Synthese künstlicher Daten auf Basis der echten Daten dar. Diese weisen idealerweise die relevanten Eigenschaften der ursprünglichen persönlichen Daten auf, enthalten aber keine tatsächlichen persönlichen Informationen.

Realistische Daten werden oft als Testdaten für Forschung und Entwicklung benötigt. Obwohl die Daten in solchen Fällen nicht real sein müssen, müssen sie realistische Eigenschaften aufweisen. In einigen dieser Anwendungsfälle können anonymisierte Daten eine Option sein. In vielen Szenarien müssen die Daten jedoch einem bestimmten Format folgen, das anonymisierte Daten aufgrund von Auslassungen und Verallgemeinerungen (je nach Anonymisierungsstrategie) oft nicht liefern können.

Eine Alternative zur Anonymisierung echter Daten könnte die Synthese von Daten sein, die die Merkmale der gewünschten Art von personenbezogenen Daten aufweisen, ohne tatsächliche persönliche Informationen zu enthalten. Dies kann durch hinreichend zufällige Variationen bestehender, wahrer personenbezogener Daten oder durch die Simulation rein synthetischer Daten auf der Grundlage geeigneter Regeln und

Annahmen geschehen.

Rein synthetische Daten basieren auf heuristischen Annahmen oder (eindeutig anonymen) Statistiken. Beispielsweise kann man Annahmen über Wahrscheinlichkeitsverteilungen, ausgearbeitete Listen mit üblichen Vor- und Nachnamen und Bevölkerungsstatistiken verwenden. Rein synthetische Daten beziehen sich garantiert nicht auf reale Personen – auch nicht bei zufälligen Kollisionen von Namen oder anderen Attributen. Der Nutzen von rein synthetischen Daten hängt von dem Zweck ab, für den die Daten verwendet werden sollen, sowie von der Qualität des Simulationsmodells und den Annahmen und Entscheidungen für die Simulation. In einigen Fällen können sogar beliebige Zufallsdaten nach der angegebenen Syntax ausreichend sein, aber in vielen anderen Fällen müssen die synthetischen Daten semantisch realistisch sein.

Die andere Klasse von synthetischen Daten sind „hybride“ synthetische Daten, die auf einem Satz von ursprünglich personenbezogenen Daten basieren. Dies entspricht der Anonymisierung durch Randomisierung (s. Kapitel 14). Die Methoden reichen von der zufälligen „Maskierung“ (d. h. Ersetzung) bis hin zur Anwendung des exponentiellen Mechanismus

<sup>1</sup> <https://css.csail.mit.edu/cryptdb/>.

<sup>2</sup> <https://github.com/google/encrypted-bigquery-client>.

für Differential Privacy (s. Unterabschnitt 14.2.2) oder von maschinellem Lernen in Verbindung mit Differential Privacy (s. Unterabschnitt 16.2.1). Solche Ansätze finden bisher in der Anwendung jedoch noch keine breite Akzeptanz. Hybride synthetische Daten können zudem Informationen über Personen enthalten, da sie noch zu viele Überreste des ursprüng-

lichen Datensatzes personenbezogener Daten enthalten können. Es gibt also einen Kompromiss zwischen Realitätsnähe und Datenschutzrisiken. Eine wesentliche Herausforderung besteht darin, Methoden zu finden, die innerhalb der Grenzen dieses Kompromisses nahezu optimale Ergebnisse erzielen und gleichzeitig recheneffizient sind.

## 17.2 Zusammenspiel von gesetzlichen Anforderungen und Technologie

Aus Sicht von Technikern erscheinen die gesetzlichen Datenschutzerfordernungen oft sehr anspruchsvoll und gleichzeitig sehr unscharf. Einige dieser Problempunkte werden im Fol-

genden angesprochen, und wir zeigen, wie das Zusammenspiel von Recht und Technik gestärkt werden muss.

### 17.2.1 Konkretisierung erforderlicher Maßnahmen und Handlungen

Die DSGVO ist sehr vage, wie die definierten Anforderungen umgesetzt werden sollen. So wird beispielsweise nicht festgelegt, welcher Aufwand für die Unterrichtung der betroffenen Personen, die die Daten nicht direkt übermittelt haben (Artikel 14), oder welcher Aufwand für die Unterrichtung der für die Löschung verantwortlichen Dritten als angemessen erachtet wird (Artikel 17). Es ist auch unklar, welche technischen und organisatorischen Maßnahmen als ausreichend oder angemessen zur Umsetzung des Datenschutzes durch Technikgestaltung (Artikel 25) und zur Gewährleistung der Sicherheit der Verarbeitung (Artikel 32) angesehen werden.

Natürlich kann das Gesetz nicht spezifischer sein, weil es alle Arten von Anwendungen abdeckt, und es muss technologie-neutral sein und auf zukünftige Entwicklungen vorbereitet sein. Daher legt die DSGVO in vielen ihrer Artikel fest, dass die jeweils nicht genauer spezifizierten Maßnahmen und Aktionen „unter Berücksichtigung“ mehrerer, spezifizierter Aspekte gewählt werden müssen. So verlangt beispielsweise Artikel 32, dass „geeignete technische und organisatorische Maßnahmen“ „[u]nter Berücksichtigung“ der folgenden vier Bedingungen gewählt werden müssen: 1. „Stand der Technik“, 2. „Implementierungskosten“, 3. „Art, Umfang, Umstände und Zwecke der Verarbeitung [Kasus angepasst]“ und 4. „Eintrittswahrscheinlichkeit und Schwere des Risikos für die Rechte und Freiheiten natürlicher Personen“. Solche Anforderungen ohne weitere Konkretisierung führen jedoch

zu einer großen Unsicherheit für die Praktiker. Insbesondere ist schwer zu sagen, ob einige bestehende, aber derzeit nicht weit verbreitete Technologien (z. B. Ansätze zur Berechnung auf verschlüsselten Daten, vgl. Kapitel 13 und Unterabschnitt 17.1.2) in bestimmten Szenarien zur Einhaltung des Gesetzes angewendet werden müssen.

Das Gesetz kann den Praktikern keine weitere Anleitung geben, da es generisch und technologie-neutral sein muss. Die DSGVO legt jedoch die Grundlage für weitere Schritte zu konkreten Anforderungen, indem sie mehrere Prozesse fördert:

1. Für die Verfeinerung der technischen und organisatorischen Maßnahmen, die in bestimmten Szenarien anzuwenden sind, empfiehlt die DSGVO die Schaffung von sektorspezifischen Verhaltensregeln (Codes of Conduct).
2. Die DSGVO fördert auch die Entwicklung von Siegeln und Prüfzeichen für die Datenschutzzertifizierung.

Einige Verbände und Einzelorganisationen haben Verhaltenskodizes entwickelt, die die gesetzlichen Anforderungen der DSGVO für ihre Branche konkretisieren. So hat die Cloud Security Alliance beispielsweise einen Verhaltenskodex<sup>3</sup> geschaffen, und der Gesamtverband der Deutschen Versicherungswirtschaft (GDV) hat seinen Verhaltenskodex<sup>4</sup> angesichts der DSGVO aktualisiert. Eine Reihe von Anbietern führen

<sup>3</sup> <https://gdpr.cloudsecurityalliance.org/>.

<sup>4</sup> <https://www.gdv.de/de/ueber-uns/unsere-services/daten-schutz-ko-dex---code-of-conduct---15544>, <https://www.gdv.de/resource/blob/23938/4aa2847df2940874559e51958a0bb350/download-code-of-conduct-data.pdf>.

Datenschutz Zertifizierungen durch und stellen entsprechende Siegel und Prüfzeichen aus. Eine Übersicht gibt es etwa bei der Stiftung Datenschutz.<sup>5</sup> Diese Angebote haben jedoch nach wie vor keine hohe Marktdurchdringung.

Eine ähnliche Entwicklung in Richtung Konkretisierung von Anforderungen vollzieht sich bei dem Thema Datenschutz-Folgenabschätzung (DSFA, engl. Data Protection Impact Assessment, DPIA). Während die DSGVO bei den Voraussetzungen und Details einer DSFA knapp gehalten ist, hat sie die Aufgabe, Details festzulegen, an die Datenschutzbehörden delegiert. Tatsächlich hatte die ehemalige Artikel-29-Arbeitsgruppe bereits vor Beginn der Anwendung der DSGVO Leitlinien für DSFAs<sup>6</sup> entwickelt. Viele nationale Behörden, auch die Datenschutzkonferenz der Datenschutzbehörden in Deutschland<sup>7</sup>, sowie private Organisationen und Verbände haben zusätzliche Spezifikationen und weitere Empfehlungen veröffentlicht.

Zusammenfassend lässt sich sagen, dass es viele Aktivitäten zur Präzisierung der gesetzlichen Anforderungen an technische und organisatorische Maßnahmen gibt. Einige dieser Richtlinien sind bereits verfügbar, und weitere sind in naher Zukunft zu erwarten. Bei den Siegeln und Prüfzeichen bleibt abzuwarten, ob sich bestimmte Angebote als Standard auf dem Markt etablieren. Möglicherweise werden zusätzliche Impulse benötigt, um etwa anbieterübergreifende Standards und Marken zu entwickeln. Institutionen wie die ENISA oder das Bundesamt für Sicherheit in der Informationstechnik (BSI) könnten die Rolle als Vertrauensanker und Koordinator für diese Entwicklung übernehmen.

---

## 17.2.2 Definition von Personenbezug und Anonymität

---

Grundsätzlich besteht eine große Schwierigkeit in einer exakten Definition von personenbezogenen Daten. Analog dazu ist es schwierig zu definieren, wann Daten nicht personenbezogen, also anonym, sind.<sup>8</sup>

Die DSGVO definiert personenbezogene Daten als solche, „die sich auf eine identifizierte oder identifizierbare natürliche Person [. . .] beziehen“ (Artikel 4 Nr. 1 DSGVO). Während die Identifizierbarkeit von Personen ein Stück weit erläutert wird, bleibt völlig un spezifiziert, wie konkret das Beziehen auf eine natürliche Person sein muss oder wie vage es sein kann, damit ein Personenbezug im Sinne des Gesetzes gegeben ist. Bei beiden Aspekten gibt es Unsicherheiten bezüglich der Auslegung, die hier erläutert werden sollen. Während der Interpretationsspielraum des Aspekts der Identifizierbarkeit in der Fachliteratur bereits umfassend erörtert ist, ist der Interpretationsspielraum beim Aspekt des Beziehens auf eine Person bisher kaum betrachtet worden.

**Wann ist eine Person identifizierbar?** Die Rechtsprechung und die Rechtswissenschaft unterscheiden zwischen der absoluten Theorie des Personenbezugs und der relativen Theorie des Personenbezugs (vgl. Kapitel 6). In der absoluten Theorie sind Daten personenbezogen, sobald die Möglichkeit besteht, dass irgendjemand auf der Welt aus diesen Daten Informationen über eine konkrete Person gewinnen kann. Die relative Theorie entscheidet basierend auf denjenigen Entitäten, die tatsächlich Zugang zu den Daten haben. Nur wenn eine dieser Instanzen die Möglichkeit hat, daraus Informationen über Personen zu gewinnen, gelten die Daten als personenbezogen im Sinne der relativen Theorie. Der Unterschied zwischen diesen Theorien hat viele Diskussionen über die Einordnung von Daten ausgelöst, insbesondere über Daten, die mit IP-Adressen verknüpft sind, wo Urteile in beide Richtungen existieren.

Darüber hinaus ist zu betrachten, wie aufwendig die Re-Identifizierung von Personen im jeweiligen Fall ist. Es gibt in bestimmten Domänen sogar gesetzliche Regelungen, die

5 [https://stiftungdatenschutz.org/fileadmin/Redaktion/PDF/Zertifizierungsuebersicht/SDS-Zertifizierungsuebersicht\\_02\\_2017.pdf](https://stiftungdatenschutz.org/fileadmin/Redaktion/PDF/Zertifizierungsuebersicht/SDS-Zertifizierungsuebersicht_02_2017.pdf).

6 [http://ec.europa.eu/newsroom/article29/item-detail.cfm?item\\_id=611236](http://ec.europa.eu/newsroom/article29/item-detail.cfm?item_id=611236).

7 [https://www.datenschutzkonferenz-online.de/media/kp/dsk\\_kpnr\\_5.pdf](https://www.datenschutzkonferenz-online.de/media/kp/dsk_kpnr_5.pdf), [https://www.datenschutzkonferenz-online.de/media/ah/20181017\\_ah\\_DSK\\_DSFA\\_Muss-Liste\\_Version\\_1.1\\_Deutsch.pdf](https://www.datenschutzkonferenz-online.de/media/ah/20181017_ah_DSK_DSFA_Muss-Liste_Version_1.1_Deutsch.pdf).

8 Vgl. die Ausführungen in Kapitel 6.

nach dem Aufwand für eine Re-Identifizierung unterscheiden und dementsprechend verschiedene Stufen von Anonymität definieren. Insbesondere unterscheidet das „Gesetz über die Statistik für Bundeszwecke“ (BStatG) zwischen formaler Anonymisierung, faktischer Anonymisierung und absoluter Anonymisierung. Es gesteht den statistischen Ämtern in Deutschland zunehmende Datenverarbeitungs- und Austauschmöglichkeiten mit zunehmender Anonymitätsstufe zu.<sup>9</sup> Aus dieser Perspektive ist es überraschend, dass die DSGVO bei der Definition personenbezogener Daten nur „schwarz“ und „weiß“ kennt.

Die Anwendung der relativen Theorie der Anonymität und von Anonymitätsstufen führt zu einer offeneren Nutzung von Daten, was für Big-Data-Anwendungen von Vorteil wäre. Der Umgang mit lediglich relativ, formal oder faktisch anonymen Daten erfordert jedoch nach wie vor angemessene Garantien und Verantwortlichkeiten unter Betrachtung der jeweiligen Datenempfänger, des Zwecks der Verarbeitung und der Risiken für den Einzelnen. Die Verarbeitung und Weitergabe solcher Daten muss kontrolliert und geregelt werden, da die Daten wieder zu personenbezogenen Daten werden können. Wir möchten betonen, dass sowohl Regeln als auch Erleichterungen für die Verwendung von Daten, die sich auf dem breiten Übergang zwischen personenbezogenen und nicht personenbezogenen Daten befinden, verfügbar sein sollten. Wir sind uns jedoch nicht sicher, wie gut der Rahmen der DSGVO geeignet ist, diesen Bereich angemessen zu regulieren. In diesem Bereich ist weitere Forschung von Rechtsexperten erforderlich.

**Wann beziehen sich Daten auf eine Person?** Bei dem Aspekt des Beziehens auf eine Person ist die alte Fassung des Bundesdatenschutzgesetzes (BDSG a. F.) etwas spezifischer als die DSGVO. Das BDSG a. F. definiert personenbezogene Daten als „Einzelangaben“ zu natürlichen Personen (§ 3 Abs. 1 BDSG a. F.). Hier wird deutlich, dass es nicht um allgemeine statistische Aussagen über Personen geht, sondern um Angaben über einzelne, konkrete Personen. Es ist jedoch zu berücksichtigen, dass in vielen Fällen auch bei Mehrpersonenangaben *Rückschlüsse* über einzelne der einbezogenen Personen getroffen werden können. Dadurch ist der Übergang zwischen Einzelangaben und anonymen Statistiken fließend, und es ist nicht klar, wo die Grenze im Sinne des BDSG a. F. liegt und noch weniger, wo sie im Sinne der DSGVO liegt.

Da in der Datenschutzrichtlinie, welche durch die DSGVO abgelöst wurde, die Definition von personenbezogenen Daten im Kern identisch mit der Definition aus der DSGVO ist, sind die Deutungen der ehemaligen Artikel-29-Datenschutzgruppe zu dem Begriff der personenbezogenen Daten auch im Kontext der DSGVO relevant. In der Opinion 4/2007 [5] untersucht und erläutert die Gruppe die gesetzliche Definition personenbezogener Daten ausführlich und geht dabei auch auf den Aspekt ein, in welcher Art und Weise sich Daten auf Personen beziehen können (d. h. auf die begriffliche Bedeutung von „relating to“). Die Möglichkeit von Rückschlüssen auf Einzelpersonen aus Informationen zu mehreren Personen wird dabei jedoch nicht behandelt. In der Opinion 05/2014 [4] befasst sich die Gruppe mit dem Thema Anonymisierung und hier werden Rückschlüsse als eines der zentralen Risiken betrachtet. Dieses Risiko wird folgendermaßen charakterisiert: „Inference, which is the possibility to deduce, with significant probability, the value of an attribute from the values of a set of other attributes“, wonach eine sehr weitreichende Definition von Rückschlüssen gewählt wird. Demnach kann der Personenbezug von Daten auch weit in allgemeine statistische Aussagen hinein bestehen.

Der Interpretationsspielraum bei dem Aspekt des Beziehens auf eine Person muss ausgeräumt werden. Entweder muss bei diesem Aspekt eine scharfe und klare Abgrenzung zwischen personenbezogenen und anonymen Daten geschaffen werden, oder es müssen analog zu dem Aspekt der Identifizierbarkeit Abstufungen geschaffen werden, die, wie bereits oben gefordert, für gesetzliche Erleichterungen beim Umgang mit Daten auf entsprechenden Stufen genutzt werden könnten. Für die Schaffung einer klaren Definition von Personenbezug und für die Schaffung möglicher Erleichterungen bei bestimmten Anonymitätsgraden ist insbesondere die Rechtsentwicklung gefragt, aber auch ein Dialog mit technischen Experten nötig. Eine präzise rechtliche Definition ist zwingende Voraussetzung für technische Anonymitätskriterien, die zweifelsfrei rechtliche Konformität garantieren können (vgl. Unterabschnitt 17.4.2).

<sup>9</sup> <https://www.forschungsdatenzentrum.de/de/anonymitaet>.

---

### 17.2.3 Verhindern von Diskriminierung

---

Scoring-Systeme auf Basis von Big Data und komplexer Algorithmen bergen ein erhöhtes Risiko für die systematische Diskriminierung bestimmter Personen in der Bevölkerung. Dies ist offensichtlich ein ethisch problematischer und unerwünschter Effekt. Darüber hinaus verbieten Gesetze und internationale Verträge eindeutig Diskriminierung, z. B. die Europäische Menschenrechtskonvention (Artikel 14), die Charta der Grundrechte der Europäischen Union (Artikel 21) und viele europäische und nationale Gesetze für spezifische Bereiche.

Um Diskriminierung durch Big Data zu verhindern, benötigen Verarbeitungs-Verantwortliche und Betreiber der entspre-

chenden Systeme geeignete organisatorische und technische Schutzmaßnahmen. Ein wichtiger Faktor kann die Beratung durch eine Ethikkommission während des gesamten Lebenszyklus der Anwendung sein, beginnend mit der Entwurfsphase bis hin zur operativen Nutzung des Systems. Weitere Maßnahmen sind eine systematische und gründliche Bewertung des Systems und Transparenzmaßnahmen. Das junge und lebendige Forschungsgebiet der algorithmischen Rechenschaftspflicht (engl. algorithmic accountability) ist ein Beweis für die Relevanz des Diskriminierungsrisikos und eine Quelle von Leitfäden und Forschungsergebnissen zu diesem interdisziplinären und komplexen Thema.

## 17.3 Unterstützung der involvierten Parteien

Im Folgenden schlagen wir Maßnahmen vor, die sich an die verschiedenen Akteure und Parteien richten, die mit Big-Data-Systemen zu tun haben.

---

### 17.3.1 Anwendbarkeit von Privacy by Design

---

Viele Prinzipien und Konzepte im Kontext von Privacy by Design scheinen einen großen Zusatzaufwand mit sich zu bringen. Für einige gibt es noch nicht einmal konkrete Ansätze zur Lösung von Praxisproblemen. Um die Idee von Privacy by Design so schnell wie möglich fest im Designprozess zu verankern, bedarf es weiterer Forschung, um Konzepte, welche die Privatheit zu akzeptablen Kosten erhöhen, klar zu identifizieren und um diese in vorkonfigurierte Systeme zu integrieren. Best Practices wären ebenfalls sehr hilfreich, um die Idee von Privacy by Design weiter zu etablieren. Ein Anfang waren die

von Hoepman [40] im Jahr 2013 vorgeschlagenen Datenschutzstrategien, die später von der ENISA [19] übernommen wurden. Die Herausforderung besteht nun darin, diese Strategien weiter zu konkretisieren und sie an die tatsächlichen Bedürfnisse im Hinblick auf die heutigen Anforderungen anzupassen. Deshalb müssen wir uns auf die Perspektive der Systembetreiber konzentrieren und konkrete Vorschläge unterbreiten. Nicht nur darüber, wie die Datenschutztechnologie leicht nutzbar gemacht werden kann, sondern in einem ersten Schritt auch darüber, wie sie eingerichtet werden kann.

---

### 17.3.2 Schulung von Praktikern

---

Eine wesentliche Aufgabe ist es, den Anwendern Kenntnisse über die entsprechenden gesetzlichen Anforderungen und deren Umsetzung in die Praxis zu vermitteln. Dazu schlagen wir vor, Aus- und Weiterbildungsmaßnahmen zu den gesetzlichen Anforderungen für die jeweilige Zielgruppe einzurichten. Ein nächster Schritt wäre es, Leitlinien und Empfehlungen vorzuschlagen, wie mit diesem gewonnenen Wissen umzugehen ist, d. h. wie die gesetzlichen Anforderungen in der

Praxis umgesetzt werden können. Sowohl für diejenigen, die Big-Data-Systeme aufsetzen, als auch für diejenigen, die diese Systeme einsetzen, empfehlen wir, explizite Beispiele mit konkreten Angreifermodellen (z. B. Bedrohungsszenarien, Art der Angreifer, Ziele und Fähigkeiten, mögliches Hintergrundwissen) durchzudenken, um geeignete Schutzmaßnahmen beurteilen und auswählen zu können.

---

### 17.3.3 Finanzielle Förderung von Forschung und Entwicklung

---

Für die Anwender von Big-Data-Technologien ist es nicht lukrativ, in die Erforschung von Datenschutztechnologien zu investieren, da dies zu einem Teufelskreis führen kann. Sie würden in Forschung investieren, die zu Lösungen führt, die den Datenschutz erhöhen und die – sobald sie einmal obligatorisch, da Stand der Technik sind – auch die Kosten der Produktivsysteme erhöhen, z. B. durch Implementierungskosten

oder Mehraufwand durch zusätzliche Berechnungen und Datenmengen. Dies führt dazu, dass es aus Anwendersicht schwer ist, für Investitionen in die Forschung zum Schutz der Privatsphäre zu argumentieren. Daher müssen kontinuierlich öffentliche Mittel zur Verfügung stehen, um die zukünftige Entwicklung in diesem Bereich zu gewährleisten.

---

### 17.3.4 Erleichterung des Zugangs zum Datenschutz für Endnutzer

---

Obwohl offensichtlich, ist es immer noch ein Punkt, wo wir Verbesserungen brauchen: Einbeziehung der Endnutzer und Erhöhung ihrer Aufmerksamkeit bzgl. des Datenschutzes und ihrer Bereitschaft, sich an der Umsetzung zu beteiligen. Außerdem müssen wir das Bewusstsein der Endnutzer schärfen und ihnen beibringen, dass sie tatsächlich etwas zum Schutz ihrer Privatsphäre tun können, wobei sie dafür keine Technologieexperten dafür sein müssen.

Eine Erleichterung für Endnutzer könnte die Einführung und Anwendung von Standard-Datenschutzrichtlinien sein, die an einem bestimmten Label leicht erkennbar sind. Beispiele für etablierte Labels finden sich in anderen Themenbereichen; hier sind insbesondere Standardlizenzen wie GPL<sup>10</sup> für

Software oder Creative Commons<sup>11</sup> für urheberrechtlich geschützte Werke zu nennen. Datenschutzzertifizierungen und entsprechende Siegel oder Prüfzeichen (vgl. Unterabschnitt 17.2.1) können ebenso datenschutzorientiertes Handeln von Endanwendern erleichtern. Nicht zu vernachlässigen, wenn auch umstritten, sind Ansätze, die spielerische Anreize (engl. gamification) zum Schutz der Privatsphäre einsetzen. Ein Beispiel könnte sein, einen guten „Privacy Score“ in Alltagsanwendungen wie Chatprogrammen zu erreichen. Nicht zuletzt brauchen wir Verbesserungen in Bezug auf Funktionalität und Einfachheit von Datenschutz (bzw. „trotz“ des Datenschutzes), um die Akzeptanz zu erhöhen und eine größere Gruppe von Endnutzern zu erreichen, d. h. nicht nur Technologiebegeisterte und Datenschutz-Evangelisten.

## 17.4 Langfristige Herausforderungen

In diesem Abschnitt werfen wir einen Blick auf Aspekte, die sich in kurzer Zeit nicht ändern werden. Einige dieser Aspekte sind schwierige wissenschaftliche Herausforderungen. Die Lösung dieser Herausforderungen wird idealerweise im nächsten Jahrzehnt oder zumindest innerhalb weniger Jahrzehnte zu

praktischen Lösungen führen. Andere Aspekte in diesem Abschnitt befassen sich mit der Entwicklung der Gesellschaft, die in vielleicht 20 Jahren einen neuen Blick auf das Thema oder neue Maßnahmen erfordern wird.

---

### 17.4.1 Effiziente voll-homomorphe Verschlüsselung

---

Homomorphe Verschlüsselung ermöglicht Berechnungen in der verschlüsselten Domäne, s. Abschnitt 13.1. Aktuelle homomorphe Verschlüsselungsverfahren erzeugen immer noch erhebliche Kosten in Bezug auf Datengröße und Rechen-

leistung. Um voll-homomorphe Verschlüsselung in die Praxis anwenden zu können, sind erhebliche Verbesserungen hinsichtlich der Effizienz solcher Systeme erforderlich.

---

<sup>10</sup> <https://www.gnu.org/licenses/gpl-3.0.en.html>.

<sup>11</sup> <https://creativecommons.org/>.

Darüber hinaus müssen mehrere Instrumente zur Vereinfachung der Anwendung von homomorpher Verschlüsselung entwickelt werden, damit die Verarbeitung homomorph verschlüsselter Daten genauso einfach programmiert werden kann, wie die Verarbeitung herkömmlicher Daten. Erstens muss es ein Werkzeug geben, das Rechenautomaten simuliert, die im Geheimtextraum eines homomorphen Verschlüsselungsverfahrens arbeiten, um das Paradigma zur Modellierung der Berechnungen von der wenig praxisfreundlichen Auswertung boolescher Schaltungen bzw. Funktionen hin zur Auswertung von Algorithmen in Rechenautomaten zu ändern. Als nächstes müssen Assembler und Compiler für

die Übersetzung von Hochsprachen in auf solchen Automaten ausführbare Programme vorhanden sein.

Zusätzlich sollte es langfristig spezielle Hardware (Prozessoren oder FPGAs) geben, die die elementaren Operationen im Geheimtextraum sowie das Bootstrapping direkt als Hardware-schaltungen implementiert, so wie heute ganz selbstverständlich beispielsweise Fließkommaoperationen in Hardware-schaltungen implementiert sind. Dies erfordert zunächst ein geeignetes, relativ effizientes homomorphes Schema. Die Ausführung auf passender Hardware könnte es dann auf praxistaugliche Effizienz bringen.

---

### 17.4.2 Beweisbare Anonymisierung

---

Neben einer präzisen formalen Definition von personenbezogenen Daten (vgl. 17.2.2) fehlt es auch an Kriterien, mit denen Daten zweifelsfrei überprüft werden können, ob ein Personenbezug vorhanden ist oder ob die Daten anonym sind. Mangels einer solchen Überprüfbarkeit gibt es keine Garantie, dass ein nach dem Stand der Technik anonymisierter Datenbestand auch tatsächlich anonym ist.

Derzeit gibt es verschiedene Maße für den Grad der Anonymität von Datentabellen als formale Kriterien für Anonymisierungsverfahren, vgl. Kapitel 14. Diese Maßnahmen drücken jedoch typischerweise eine aus technischer Sicht intuitive Vorstellung von Anonymität aus, die sich aus einer bestimmten Art von Angriff und bestimmten Annahmen über die Daten ergibt. Daher erwiesen sich die meisten Maße in anderen Szenarien, die von den ursprünglich betrachteten Angriffen abweichen, als unzureichend. Für unstrukturierte Daten fehlen Anonymitätskriterien sogar gänzlich.

Neue Anonymisierungskonzepte müssen auf der Grundlage klarer und formaler Definitionen von Angreifermodellen unter Angabe von Wissen, Fähigkeiten und Zielen der Angreifer entwickelt werden. Nur dann werden die Garantien und Grenzen dieser Konzepte von Anfang an sichtbar sein und die Fehler der Vergangenheit vermieden.

Irgendwann in der Zukunft könnte die Forschung eine ultimative Definition von Anonymität erreichen, die nachweisbar Angriffen in jedem Angreifermodell standhält. Diese Definition sollte vorzugsweise ein Kriterium sein, das tatsächlich an Daten überprüft werden kann. Ein solches Kriterium führt zu einer formalen Unterscheidung zwischen personenbezogenen und nicht personenbezogenen Daten. Darüber hinaus sollte ein effizienter Algorithmus zur Anonymisierung von Datenbeständen gemäß dieser Definition gefunden werden. Der Algorithmus sollte so viele Informationen wie möglich behalten, um den Nutzen des anonymisierten Datensatzes zu maximieren.

Differential Privacy geht in die richtige Richtung, und es könnte sogar nahe an der endgültigen Definition von Anonymität liegen. In Bezug auf Differential Privacy wünschen wir uns jedoch ein allgemeines Kriterium für die Wahl des richtigen Wertes für den Parameter Epsilon. Darüber hinaus stellt sich die Frage, ob es eine ebenso starke Datenschutzgarantie gibt, die zu weniger Informationsverlust führt als Differential Privacy. Grundsätzlich ist zu klären, wie viel Informationsverlust zum Erreichen von Anonymität zwingend notwendig ist, und welches Kriterium keine Anforderungen stellt, die dieses Mindestmaß an Informationsverlust künstlich erhöhen. Daher ist der Wettbewerb um das beste Anonymitätskriterium und die beste Anonymisierungsmethode noch offen.

---

### 17.4.3 Rechtliche Harmonisierung

---

Innerhalb der EU, aber auch darüber hinaus, ist die DSGVO ein großer Schritt in die richtige Richtung zur Rechtsharmonisierung. Wegen ihrer Schwächen sowie aufgrund der Weiterentwicklung von Technologien brauchen wir jedoch wahrscheinlich eine nächste Version – DSGVO 2.0 – in vielleicht 20 Jahren. In diesem Zusammenhang müssen wir abwarten, welche Folgen die DSGVO für die ganze Welt haben wird, d.h. wir müssen ihre Auswirkungen auf globale Dienstleister wie Facebook bzw. auf die zu dieser Zeit dann relevanten Dienstleister untersuchen. Zuvor müssen wir beobachten, wie die nationale Gesetzgebung in der EU und die Dienstleister auf dem Markt die von der DSGVO auf nationaler Ebene

zugelassenen Spielräume nutzen. Die DSGVO 2.0 kann dann basierend auf diesen Beobachtungen eine sinnvolle weitere Harmonisierung in der EU erreichen.

Im Idealfall wird es in der EU und darüber hinaus zu einer marktgetriebenen Konvergenz von Datenschutzstandards kommen. Andernfalls wird der Weg zu einer (globalen) Rechtsharmonisierung deutlich schwieriger und schleppender sein, auch wenn der California Consumer Privacy Act als Ausstrahlungserfolg der DSGVO gewertet werden kann und ein erster Schritt in Richtung Annäherung von Datenschutzrechten ist.

---

### 17.4.4 Nicht absehbare Entwicklungen der Gesellschaft

---

Nicht zuletzt müssen wir uns bei einem Blick in die fernere Zukunft mit noch unbekanntem Faktoren der zukünftigen Gesellschaft auseinandersetzen. Ein Beispiel von vor wenigen Jahrzehnten ist die Volkszählung in Deutschland in den 1980er Jahren. Die massiven Proteste gegen die Erhebung führten dazu, dass diese zunächst nicht durchgeführt wurde. Die Gesellschaft war nicht bereit, der Regierung eine solche Menge an privaten Daten zu geben. Die Bürger fühlten sich im Recht auf Privatsphäre stark verletzt und das Bundesverfassungsgericht leitete im als Volkszählungsurteil in die Geschichte eingegangenen Urteil das Recht auf informelle Selbstbestimmung aus den verfassungsmäßig zugesicherten Grundrechten ab. Im Gegensatz dazu sind heute viele Personen bereit, viel mehr Informationen nicht nur an Regierungsbehörden, sondern auch an große Technologieunternehmen oder sogar an die Öffentlichkeit weiterzugeben, indem sie private Informationen in sozialen Netzwerken veröffentlichen oder Smartphones benutzen.

Ein weiterer Faktor, der die Unsicherheit über zukünftige Datenschutznormen erhöht, ist die mögliche Verschiebung der politischen Rahmenbedingungen in Abhängigkeit von zukünftigen Sicherheitsbedürfnissen und wirtschaftlichen Entwicklungen. So könnte die Schwächung der Privatsphäre zugunsten von wirtschaftlichem Nutzen oder für Strafverfolgung und Überwachung aus politischen Interessen durchgesetzt werden. Ob dies in der jeweiligen Situation eher gesellschaftlich akzeptiert oder kritisiert werden wird, ist offen.

Diese unbekanntem Entwicklungen ist nicht wie ein technisches Problem durch das Finden von Lösungen anzugehen und in den Griff zu bekommen, aber sie sind Bedingungen, die die Erforschung und Entwicklung des Datenschutzes langfristig herausfordern könnten.

Daraus lernen wir, dass wir nicht vorhersagen können, was heute unvorstellbar ist, aber in Zukunft ganz normal. Daher müssen wir uns bewusst sein, dass Änderungen der allgegenwärtigen Sichtweisen und Verhaltensnormen die Entwicklung des Datenschutzes beeinflussen werden.



## 17.5 Fazit

Der Überblick, den wir in diesem Kapitel dargestellt haben, beschreibt die Probleme, die in Zukunft gelöst werden müssen, um Big Data datenschutzverträglicher zu machen. Wir haben versucht zu identifizieren, was in den nächsten Jahren vorangetrieben und erreicht werden sollte und kann, und was ein langfristiges Ziel für die nächsten ein oder mehreren Jahrzehnte ist.

Wir haben konkret in folgenden Bereichen zu verbessernde Punkte aufgezeigt: Bei verschiedenen Technologien, insbesondere bei der Anonymisierung, beim Rechnen auf verschlüssel-

ten Daten und bei der Synthese von Daten gibt es eine Vielzahl ungelöster Probleme, die im Laufe der Zeit durch Forschung gelöst werden müssen. Die DSGVO bringt viel Unsicherheit in das Feld, das mit entsprechenden Vorgaben und Leitfäden gefüllt werden muss. Das Wissen über Technologien und vorbildlichen Praxisumsetzungen zur Erfüllung der gesetzlichen Anforderungen muss auf diejenigen übertragen werden, die Big-Data-Systeme aufbauen und betreiben. Schließlich müssen die Endnutzer aufgeklärt werden, ihnen müssen einfach nutzbare Informationen und Werkzeuge gegeben werden und sie müssen ihre Entscheidungen treffen.

## 18 Publikationen

In diesem Projekt haben die Projektmitarbeiter folgende **Publikationen** erstellt:

1. Marcel Schäfer, Jamal Pasha, Martin Steinebach: Introducing a Hybrid Generalization and Micro-Aggregation Algorithm for Database Anonymization to reduce the Gap between Data Utility and Privacy Requirements. Amsterdam Privacy Conference 2018, 5–8 October 2018, Amsterdam, the Netherlands. Kein veröffentlichter Konferenzband.
2. Christian Winter, Marcel Schäfer: Roadmap to Privacy. Amsterdam Privacy Conference 2018, 5–8 October 2018, Amsterdam, the Netherlands. Kein veröffentlichter Konferenzband.
3. Christian Winter, Verena Battis, Oren Halvani: Herausforderungen für die Anonymisierung von Daten. Informatik 2019, 49. GI-Jahrestagung, 23.–26. September 2019, Kassel, Deutschland. Lecture Notes in Informatics, Volume P-294, Seiten 339–352. Gesellschaft für Informatik e.V., Bonn, September 2019.
4. Christian Winter, Verena Battis, Oren Halvani: Herausforderungen für die Anonymisierung von Daten. Zeitschrift für Datenschutz (ZD), Ausgabe 11/2019, Seiten 489–493. Verlag C.H.Beck oHG, München, November 2019.

Im Rahmen des Projekts haben die Projektmitarbeiter folgende **Abschlussarbeit** betreut:

1. Antonio Odoguardi: Konzeption und Implementierung eines Mikroaggregationsalgorithmus zur Datenanonymisierung auf Basis des Maximum-Distance-to-Average-Vector (MDAV)-Prinzips. Masterarbeit, TU Darmstadt, November 2019.

Zudem wurden außerhalb des Projekts folgende **Abschlussarbeiten** erstellt, die aufgrund ihres Projektbezugs in Teil III referenziert wurden und daher der Vollständigkeit halber noch einmal hier aufgelistet werden:

1. Jamal Pasha: Database Anonymization: Data Utility and Privacy Leakage in Generalization and Microaggregation. Master Thesis, TU Darmstadt, März 2017.
2. Roey Regev: Verwendung von Koreferenz-Auflösung zur Anonymisierung von Eigennamen in deutschsprachigen Texten. Bachelorarbeit, TU Darmstadt, August 2019.

## Literatur

- [1] Martín Abadi und Joan Feigenbaum. »Secure circuit evaluation – A protocol based on hiding information from an oracle«. In: *Journal of Cryptology* 2.1 (Feb. 1990), S. 1–12. ISSN: 0933-2790 (print) and 1432-1378 (online). DOI: 10.1007/BF02252866.
- [2] Nadhem J. AlFardan und Kenneth G. Paterson. »Lucky thirteen: Breaking the TLS and DTLS record protocols«. In: *2013 IEEE Symposium on Security and Privacy*. IEEE, 2013, S. 526–540.
- [3] Nadhem J. AlFardan u. a. »On the security of RC4 in TLS«. In: *22nd USENIX Security Symposium*. 2013, S. 305–320. URL: <http://www.isg.rhul.ac.uk/tls/RC4biases.pdf>.
- [4] Artikel-29-Datenschutzgruppe. *Opinion 05/2014 on Anonymisation Techniques*. Techn. Ber. WP216. Artikel-29-Datenschutzgruppe, 10. Apr. 2014. URL: [https://ec.europa.eu/justice/article-29/documentation/opinion-recommendation/files/2014/wp216\\_en.pdf](https://ec.europa.eu/justice/article-29/documentation/opinion-recommendation/files/2014/wp216_en.pdf).
- [5] Artikel-29-Datenschutzgruppe. *Opinion 4/2007 on the concept of personal data*. Techn. Ber. WP 136. Artikel-29-Datenschutzgruppe, 20. Juni 2007. URL: [https://ec.europa.eu/justice/article-29/documentation/opinion-recommendation/files/2007/wp136\\_en.pdf](https://ec.europa.eu/justice/article-29/documentation/opinion-recommendation/files/2007/wp136_en.pdf).
- [6] Nimrod Aviram u. a. »DROWN: Breaking TLS using SSLv2«. In: *25th USENIX Security Symposium*. Austin, TX: USENIX Association, 2016, S. 689–706. ISBN: 978-1-931971-32-4. URL: <https://www.usenix.org/conference/usenixsecurity16/technical-sessions/presentation/aviram>.
- [7] Josh Benaloh. »Dense probabilistic encryption«. In: *Proceedings of the workshop on selected areas of cryptography*. 1994, S. 120–128.
- [8] Karthikeyan Bhargavan und Gaëtan Leurent. »On the practical (in-)security of 64-bit block ciphers: Collision attacks on HTTP over TLS and OpenVPN«. In: *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2016, S. 456–467.
- [9] Avrim Blum, Katrina Ligett und Aaron Roth. »A Learning Theory Approach to Non-Interactive Database Privacy«. In: *Proceedings of the 40th annual ACM Symposium on Theory of Computing*. STOC’08. ACM, 2008, S. 609–618.
- [10] Keith Bonawitz u. a. »Practical Secure Aggregation for Privacy-Preserving Machine Learning«. In: *CCS 2017*. ACM, 2017, S. 1175–1191.
- [11] Dan Boneh, Eu-Jin Goh und Kobbi Nissim. »Evaluating 2-DNF Formulas on Ciphertexts«. In: *Theory of Cryptography. Second Theory of Cryptography Conference, TCC 2005, Cambridge, MA, USA, February 10–12, 2005. Proceedings*. Hrsg. von Joe Kilian. 2005, S. 325–341.
- [12] Joan Boyar, René Peralta und Denis Pochuev. »On the Multiplicative Complexity of Boolean Functions over the Basis  $(\wedge, \oplus, 1)$ «. In: *Theoretical Computer Science* 235.1 (März 2000), S. 43–57.
- [13] Gilles Brassard und Claude Crépeau. »Zero-knowledge Simulation of Boolean Circuits«. In: *Advances in Cryptology – CRYPTO ’86*. Hrsg. von Andrew M. Odlyzko. Lecture Notes in Computer Science (LNCS) 263. Springer, 1987, S. 223–233. ISBN: 978-3-540-18047-0 (print) and 978-3-540-47721-1 (online). DOI: 10.1007/3-540-47721-7\_16.
- [14] Bundesamt für Sicherheit in der Informationstechnik (BSI). *Kryptographische Verfahren: Empfehlungen und Schlüssellängen. Teil 2 – Verwendung von Transport Layer Security (TLS)*. Technische Richtlinie TR-02102-2. BSI, 2016. URL: [https://www.bsi.bund.de/SharedDocs/Downloads/DE/BSI/Publikationen/TechnischeRichtlinien/TR02102/BSI-TR-02102-2.pdf?\\_\\_blob=publicationFile&v=2](https://www.bsi.bund.de/SharedDocs/Downloads/DE/BSI/Publikationen/TechnischeRichtlinien/TR02102/BSI-TR-02102-2.pdf?__blob=publicationFile&v=2).
- [15] Bundesgerichtshof. *Urteil*. Aktenzeichen VI ZR 156/13. Jan. 2014. URL: <http://juris.bundesgerichtshof.de/cgi-bin/rechtsprechung/document.py?Gericht=bgh&nr=66910>.

- [16] Ann Cavoukian. *Privacy by Design: The 7 Foundational Principles*. Techn. Ber. Revised in January 2011. Information und Privacy Commissioner, Ontario, Canada, Aug. 2009. URL: <https://www.ipc.on.ca/wp-content/uploads/Resources/7foundationalprinciples.pdf>.
- [17] Ann Cavoukian und Jeff Jonas. *Privacy by Design in the Age of Big Data*. Techn. Ber. Information und Privacy Commissioner, Ontario, Canada, Juni 2012. URL: <https://jeffjonas.typepad.com/Privacy-by-Design-in-the-Era-of-Big-Data.pdf>.
- [18] H. Cheng und Xiaobo Li. »Partial Encryption of Compressed Images and Videos«. In: *IEEE Transactions on Signal Processing* 48.8 (Aug. 2000), S. 2439–2451.
- [19] George Danezis u. a. *Privacy and Data Protection by Design – from policy to engineering*. ENISA, Dez. 2014. ISBN: 978-92-9204-108-3. URL: <https://www.enisa.europa.eu/activities/identity-and-trust/library/deliverables/privacy-and-data-protection-by-design>.
- [20] Marten van Dijk u. a. »Fully Homomorphic Encryption over the Integers«. In: *Advances in Cryptology – EUROCRYPT 2010*. Bd. 6110. Lecture Notes in Computer Science. Springer, 2010, S. 24–43.
- [21] Inga Döbel u. a. *Maschinelles Lernen – Eine Analyse zu Kompetenzen, Forschung und Anwendung*. Report. Fraunhofer-Gesellschaft, 2018.
- [22] Josep Domingo-Ferrer und Josep M. Mateo-Sanz. »Practical Data-Oriented Microaggregation for Statistical Disclosure Control«. In: 14.1 (Jan. 2002), S. 189–202. DOI: 10.1109/69.979982.
- [23] Nathan Dowlin u. a. *CryptoNets: Applying Neural Networks to Encrypted Data with High Throughput and Accuracy*. Techn. Ber. 2016. URL: <https://www.microsoft.com/en-us/research/publication/cryptonets-applying-neural-networks-to-encrypted-data-with-high-throughput-and-accuracy/>.
- [24] Yang Du u. a. »On Multidimensional k-Anonymity with Local Recoding Generalization«. In: *Proceeding of 2007 IEEE 23rd International Conference on Data Engineering*. ICDE 2007. IEEE Computer Society, 2007, S. 1422–1424. DOI: 10.1109/ICDE.2007.369026.
- [25] John C. Duchi, Michael I. Jordan und Martin J. Wainwright. »Local Privacy and Statistical Minimax Rates«. In: *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*. FOCS 2013. IEEE Computer Society Conference Publishing Service (CPS), 2013, S. 429–438. DOI: 10.1109/FOCS.2013.53.
- [26] Thai Duong und Juliano Rizzo. *Here Come The  $\oplus$  Ninjas*. 2011. URL: <http://www.hpcc.ecs.soton.ac.uk/~dan/talks/bullrun/Beast.pdf>.
- [27] Cynthia Dwork. »Differential Privacy«. In: *Automata, Languages and Programming*. Springer, 2006, S. 1–12.
- [28] Cynthia Dwork. »Differential Privacy: A Survey of Results«. In: *Theory and Applications of Models of Computation. 5th International Conference, TAMC 2008, Xi'an, China, April 25–29, 2008. Proceedings*. Hrsg. von Manindra Agrawal u. a. Bd. 4978. Lecture Notes in Computer Science. Springer, 2008, S. 1–19.
- [29] Cynthia Dwork u. a. »Calibrating Noise to Sensitivity in Private Data Analysis«. In: *Theory of Cryptography*. Springer, 2006, S. 265–284.
- [30] Cynthia Dwork u. a. »Calibrating noise to sensitivity in private data analysis«. In: *Journal of Privacy and Confidentiality* 7.3 (2016), S. 17–51.
- [31] Morris Dworkin. *Recommendation for Block Cipher Modes of Operation: The XTS-AES Mode for Confidentiality on Storage Devices*. Special Publication (SP) 800-38E. National Institute of Standards und Technology (NIST), Jan. 2010. URL: <https://csrc.nist.gov/publications/detail/sp/800-38e/final>.
- [32] Europarat. *Konvention zum Schutz der Menschenrechte und Grundfreiheiten (Ursprüngliche Fassung auf Englisch und Französisch)*. SEV Nr.005. Europarat, 4. Nov. 1950. URL: <https://www.coe.int/de/web/conventions/full-list/-/conventions/treaty/005>. In Kraft seit 3. September 1953, zuletzt geändert durch Zusatzprotokoll Nr. 14 (SEV Nr.214) vom 2. Oktober 2013, welches seit 1. August 2018 in Kraft ist.

- [33] Matt Fredrikson, Somesh Jha und Thomas Ristenpart. »Model inversion attacks that exploit confidence information and basic countermeasures«. In: *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*. ACM CCS 2015. ACM, 2015, S. 1322–1333.
- [34] Craig Gentry. »A fully homomorphic encryption scheme«. Diss. Stanford University, Sep. 2009. URL: <https://crypto.stanford.edu/craig/>.
- [35] Craig Gentry. »Fully homomorphic encryption using ideal lattices«. In: *Proceedings of the forty-first annual ACM symposium on Theory of computing*. STOC'09. ACM, 2009, S. 169–178.
- [36] Shafi Goldwasser und Silvio Micali. »Probabilistic encryption«. In: *Journal of computer and system sciences* 28.2 (1984), S. 270–299.
- [37] Tommi Gröndahl und N. Asokan. *Text Analysis in Adversarial Settings: Does Deception Leave a Stylistic Trace?* 26. Feb. 2019. arXiv: 1902.08939v2 [cs.CL].
- [38] H. Hakami und D. Bollegala. »Learning Relation Representations from Word Representations«. In: *Submitted to Automated Knowledge Base Construction. under review*. 2019. URL: <https://openreview.net/forum?id=r1e3WW5aTX>.
- [39] Briland Hitaj, Giuseppe Ateniese und Fernando Pérez-Cruz. »Deep models under the GAN: information leakage from collaborative deep learning«. In: *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. ACM CCS 2017. ACM, 2017, S. 603–618.
- [40] Jaap-Henk Hoepman. *Privacy Design Strategies*. 6. Mai 2013. arXiv: 1210.6621v2 [cs.CY].
- [41] Shiva Prasad Kasiviswanathan u. a. »What Can We Learn Privately?«. In: *49th Annual IEEE Symposium on Foundations of Computer Science*. FOCS 2008. IEEE Computer Society, 2008, S. 531–540.
- [42] Thomas Kunkelmann. »Partielle Verschlüsselung«. In: *Sicherheit für Videodaten*. Vieweg+Teubner Verlag, 1998, S. 105–124. DOI: 10.1007/978-3-663-05777-2\_5.
- [43] Kristin Lauter, Michael Naehrig und Vinod Vaikuntanathan. »Can Homomorphic Encryption be Practical?«. In: *Proceedings of the 3rd ACM Cloud Computing Security Workshop*. CCSW'11. ACM, 2011, S. 113–124.
- [44] Haoran Li, Li Xiong und Xiaoqian Jiang. »Differentially private synthesization of multi-dimensional data using copula functions«. In: *Advances in Database Technology — EDBT 2014*. EDBT 2014. OpenProceedings, 2014, S. 475–486. URL: <http://openproceedings.org/2014/conf/edbt/LiXJ14.pdf>.
- [45] Jing Li u. a. *A Survey on Deep Learning for Named Entity Recognition*. 22. Dez. 2018. arXiv: 1812.09449v1 [cs.CL].
- [46] Ninghui Li, Tiancheng Li und Suresh Venkatasubramanian. »t-Closeness: Privacy Beyond k-Anonymity and l-Diversity«. In: *2007 IEEE 23rd International Conference on Data Engineering*. ICDE 2007. IEEE Computer Society, Apr. 2007, S. 106–115.
- [47] Yehuda Lindell und Benny Pinkas. »Privacy Preserving Data Mining«. In: *Advances in Cryptology – CRYPTO 2000. 20th Annual International Cryptology Conference Santa Barbara, California, USA, August 20–24, 2000 Proceedings*. Hrsg. von Mihir Bellare. Bd. 1880. Lecture Notes in Computer Science (LNCS). Springer, 2000, S. 36–54. DOI: 10.1007/3-540-44598-6\_3.
- [48] Yehuda Lindell und Benny Pinkas. »Secure Multiparty Computation for Privacy-Preserving Data Mining«. In: *Journal of Privacy and Confidentiality* 1.1 (2009), S. 59–98. DOI: 10.29012/jpc.v1i1.566.
- [49] Quiang Liu u. a. »A survey on security threats and defensive techniques of machine learning: A data driven view«. In: *IEEE Access* 6 (2018), S. 12103–12117.
- [50] Yunhui Long u. a. *Distributed and Secure ML with Self-tallying Multi-party Aggregation*. 26. Nov. 2018. arXiv: 1811.10296v1 [cs.CR].
- [51] Ashwin Machanavajjhala u. a. »l-Diversity: Privacy

- Beyond k-Anonymity«. In: *Proceedings of the 22nd International Conference on Data Engineering*. ICDE'06. IEEE Computer Society, Apr. 2006. DOI: 10.1109/ICDE.2006.1.
- [52] Ashwin Machanavajjhala u. a. »I-Diversity: Privacy Beyond k-Anonymity«. In: *ACM Transactions on Knowledge Discovery from Data* 1.1, 3 (März 2007).
- [53] Claire McKay Bowen und Fang Liu. *Comparative study of differentially private data synthesis methods*. 8. Jan. 2019. arXiv: 1602.01063v4 [stat.ME].
- [54] H. Brendan McMahan u. a. »Communication-Efficient Learning of Deep Networks from Decentralized Data«. In: *Artificial Intelligence and Statistics*. AISTATS 2017. Hrsg. von Aarti Singh und Jerry Zhu. Bd. 54. PMLR. PMLR, 2017, S. 1273–1282. URL: <http://proceedings.mlr.press/v54/mcmahan17a.html>.
- [55] Frank McSherry und Kunal Talwar. »Mechanism Design via Differential Privacy«. In: *Proceedings of the 48th Annual IEEE Symposium on Foundations of Computer Science*. FOCS 2007. IEEE Computer Society, 2007, S. 94–103.
- [56] Luca Melis u. a. »Exploiting unintended feature leakage in collaborative learning«. In: *2019 IEEE Symposium on Security and Privacy (S&P 2019)*. IEEE S&P 2019. 2019, S. 497–512.
- [57] UN-Menschenrechtskommission. *Allgemeine Erklärung der Menschenrechte*. Resolution der UN-Generalversammlung A/RES/217(III) A. In: *Official Records of the Third Session of the General Assembly, Part I, Resolutions (A/810)*, S. 71–77. Vereinte Nationen, 10. Dez. 1948. URL: <https://www.un.org/en/documents/udhr/>.
- [58] Bodo Möller, Thai Duong und Krzysztof Kotowicz. *This POODLE Bites: Exploiting The SSL 3.0 Fallback*. 2014. URL: <https://www.openssl.org/~bodo/ssl-poodle.pdf>.
- [59] M. Ercan Nergiz, Maurizio Atzori und Christopher W. Clifton. »Hiding the Presence of Individuals from Shared Databases«. In: *Proceedings of the ACM SIGMOD International Conference on Management of Data*. SIGMOD'07. ACM, 2007, S. 665–676. DOI: 10 . 1145 / 1247480 . 1247554.
- [60] Antonio Odoguardi. »Konzeption und Implementierung eines Mikroaggregationsalgorithmus zur Datenanonymisierung auf Basis des Maximum-Distance-to-Average-Vector (MDAV)-Prinzips«. Masterarbeit. Technische Universität Darmstadt, 2019. In Bearbeitung.
- [61] Hector Page, Charlie Cabot und Kobbi Nissim. *Differential privacy: an introduction for statistical agencies. A contributing article to the National Statistician's Quality Review into Privacy and Data Confidentiality Methods*. Techn. Ber. Government Statistical Service (GSS), 13. Dez. 2018. URL: [https://gss.civilservice.gov.uk/wp-content/uploads/2018/12/12-12-18\\_FINAL\\_Privitar\\_Kobbi\\_Nissim\\_article.pdf](https://gss.civilservice.gov.uk/wp-content/uploads/2018/12/12-12-18_FINAL_Privitar_Kobbi_Nissim_article.pdf).
- [62] Jamal Pasha. »Database Anonymization: Data Utility and Privacy Leakage in Generalization and Microaggregation«. Master Thesis. Technische Universität Darmstadt, 6. März 2017.
- [63] Andrey Popov. *Prohibiting RC4 cipher suites*. RFC 7465. Internet Engineering Task Force (IETF), Feb. 2015.
- [64] Martin Potthast u. a. »A Decade of Shared Tasks in Digital Text Forensics at PAN«. In: *Advances in Information Retrieval*. Hrsg. von Leif Azzopardi u. a. Springer, 2019, S. 291–300.
- [65] Roey Regev. »Verwendung von Koreferenz-Auflösung zur Anonymisierung von Eigennamen in deutschsprachigen Texten«. Bachelorarbeit. Technische Universität Darmstadt, Aug. 2019.
- [66] Ronald L Rivest, Len Adleman und Michael L Dertouzos. »On data banks and privacy homomorphisms«. In: *Foundations of secure computation* 4.11 (1978), S. 169–180.
- [67] Ronald L Rivest, Adi Shamir und Len Adleman. »A method for obtaining digital signatures and public-key cryptosystems«. In: *Communications of the ACM* 21.2 (1978), S. 120–126.

- [68] Mohammad Al-Rubaie und J Morris Chang. »Privacy-Preserving Machine Learning: Threats and Solutions«. In: *IEEE Security & Privacy* 17.2 (2019), S. 49–58.
- [69] Pierangela Samarati und Latanya Sweeney. »Generalizing Data to Provide Anonymity when Disclosing Information (Abstract)«. In: *Proceedings of the Seventeenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems*. PODS 1998 (Seattle, Washington, USA). New York, NY, USA: ACM, 1998, S. 188. ISBN: 0-89791-996-3. DOI: 10.1145/275487.275508.
- [70] Tomas Sander, Adam L. Young und Moti Yung. »Non-Interactive CryptoComputing for NC1«. In: *40th Annual Symposium on Foundations of Computer Science*. 1999, S. 554–566.
- [71] Sawvas Savvides u. a. »Secure data types: a simple abstraction for confidentiality-preserving data analytics«. In: *Proceedings of the 2017 Symposium on Cloud Computing (STOC '17)*. 2017, S. 479–492. DOI: 10.1145/3127479.3129256.
- [72] Marcel Schäfer, Jamal Pasha und Martin Steinebach. »Introducing a Hybrid Generalization and Micro-Aggregation Algorithm for Database Anonymization to reduce the Gap between Data Utility and Privacy Requirements«. In: *Amsterdam Privacy Conference 2018* (5.–8. Okt. 2018). 2018.
- [73] A. Servetti, C. Testa und J. C. De Martin. »Frequency-Selective Partial Encryption of Compressed Audio«. In: *2003 IEEE International Conference on Acoustics, Speech, and Signal Processing, 2003. Proceedings. (ICASSP '03)*. ICASSP '03. Bd. 5. Apr. 2003, S. V-668–V-671. DOI: 10.1109/ICASSP.2003.1200059.
- [74] Rakshith Shetty, Bernt Schiele und Mario Fritz. »A<sup>4</sup>NT: Author Attribute Anonymity by Adversarial Training of Neural Machine Translation«. In: *USENIX Security '18*. USENIX Association, 2018, S. 1633–1650.
- [75] Reza Shokri und Vitaly Shmatikov. »Privacy-preserving deep learning«. In: *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*. CCS 2015. ACM, 2015, S. 1310–1321.
- [76] Reza Shokri u. a. »Membership inference attacks against machine learning models«. In: *2017 IEEE Symposium on Security and Privacy (S&P)*. IEEE S&P 2017. 2017, S. 3–18.
- [77] Ravid Shwartz-Ziv und Naftali Tishby. *Opening the Black Box of Deep Neural Networks via Information*. 29. Apr. 2017. arXiv: 1703.00810v3 [cs.LG].
- [78] Julian James Stephen u. a. »Practical Confidentiality Preserving Big Data Analysis«. In: *6th USENIX Workshop on Hot Topics in Cloud Computing*. HotCloud 14. Philadelphia, PA: USENIX Association, Juni 2014. URL: <https://www.usenix.org/conference/hotcloud14/workshop-program/presentation/stephen>.
- [79] Marshall Harvey Stone. »The Theory of Representation for Boolean Algebras«. In: *Transactions of the American Mathematical Society* 40.1 (Juli 1936), S. 37–111. DOI: 10.2307/1989664. URL: <https://www.jstor.org/stable/1989664>.
- [80] Latanya Sweeney. »k-Anonymity: A Model for Protecting Privacy«. Englisch. In: *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems* 10.5 (Okt. 2002), S. 557–570. DOI: 10.1142/S0218488502001648.
- [81] Latanya Sweeney. »Achieving k-Anonymity Privacy Protection Using Generalization and Suppression«. Englisch. In: *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems* 10.5 (Okt. 2002), S. 571–588. DOI: 10.1142/S021848850200165X.
- [82] Latanya Sweeney. »Computational Disclosure Control. A Primer on Data Privacy Protection«. Diss. Massachusetts Institute of Technology, Mai 2001. URL: <http://hdl.handle.net/1721.1/8589>.
- [83] Aleksei Triastcyn und Boi Faltings. »Generating artificial data for private deep learning«. In: *Proceedings of the PAL: Privacy-Enhancing Artificial Intelligence and Language Technologies*. PAL. Bd. Vol-2335. CEUR Workshop Proceedings. 18. März 2019, S. 33–40.

- [84] Qian Wang, Zhiwei Xu und Shengzhi Qu. »An Enhanced K-Anonymity Model against Homogeneity Attack«. In: *Journal of Software* 6.10 (2011). URL: <http://ojs.academypublisher.com/index.php/jsw/article/view/jsw061019451952>.
- [85] Stanley L. Warner. »Randomized Response: A Survey Technique for Eliminating Evasive Answer Bias«. In: *Journal of the American Statistical Association* 60.309 (1965), S. 63–69. DOI: 10.1080/01621459.1965.10480775.
- [86] Vikas Yadav und Steven Bethard. »A Survey on Recent Advances in Named Entity Recognition from Deep Learning models«. In: *Proceedings of the 27th International Conference on Computational Linguistics. COLING 2018. Association for Computational Linguistics (ACL), Aug. 2018, S. 2145–2158*.
- [87] Andrew C. Yao. »Protocols for secure computations«. In: *23rd Annual Symposium on Foundations of Computer Science. FOCS 1982 (Chicago, Illinois, USA, 3.–5. Nov. 1982)*. IEEE Computer Society, 1982, S. 160–164. DOI: 10.1109/SFCS.1982.88.
- [88] Ivan Ivanovich Zhegalkin. »On the technique of calculating propositions in symbolic logic (in Russian, with French abstract)«. In: *Matematicheskii Sbornik* 34.1 (1927), S. 9–28. URL: <http://mi.mathnet.ru/eng/msb7433>.
- [89] Ivan Ivanovich Zhegalkin. »The arithmetization of symbolic logic (in Russian, with French abstract)«. In: *Matematicheskii Sbornik* 35.3–4 (1928), S. 311–377. URL: <http://mi.mathnet.ru/eng/msb7400>.



