

Datenschutzziele im Konflikt: Eine Metrik für Unverkettbarkeit als Hilfestellung für den Betroffenen

Christoph Bier¹

Abstract: Die 6 Datenschutz-Schutzziele befinden sich in einem dauerhaften Spannungsverhältnis. Insbesondere Transparenz und Unverkettbarkeit kollidieren. Das Recht auf Auskunft ist ein prominentes Beispiel für den Konflikt zwischen den beiden Schutzzielen. Es erfordert eine umfangreiche Rückverfolgung und mithin eine stärkere Verkettung personenbezogener Daten um den Betroffenen über den Umgang mit denselben in Kenntnis zu setzen. Im Sinne eines mündigen Bürgers sollte der Betroffene entscheiden, was ihm wichtiger ist: Die Unverkettbarkeit oder eine möglichst umfangreiche Auskunft. Dies bedarf jedoch einer klaren Kommunikation der Konsequenzen.

In diesem Beitrag wird auf Grundlage rechtlicher Anforderungen eine informationstheoretische Metrik für Unverkettbarkeit vorgestellt, modelliert und für ein automatisiertes Datenschutzauskunftssystem instantiiert. Die Metrik bildet die notwendige Informationsgrundlage, um es Betroffenen zu erlauben, selbst über den Trade-off zwischen Transparenz und Unverkettbarkeit zu entscheiden. Eine prototypische Implementierung zeigt die Berechenbarkeit der Metrik während das Ergebnis einer Nutzerstudie mit 31 Teilnehmern eine gute Akzeptanz der Metrik nahelegt.

Keywords: Datenschutz, Unverkettbarkeit, Transparenz, Auskunft, Metrik, Entropie

1 Einleitung

Das Datenschutzrecht versucht mit einer Vielzahl von Regelungen das informationelle Selbstbestimmungsrecht der Bürger und deren Handlungs- und Mitwirkungsfähigkeit im freiheitlich-demokratischen Gemeinwesen der Bundesrepublik Deutschland sicherzustellen.² Das Bundesdatenschutzgesetz (BDSG) und in die EU-Datenschutzgrundverordnung (EU-DSGVO) enthalten unterschiedlichste Anforderungen an die Verarbeitung personenbezogener Daten. Die in den letzten Jahren entwickelten Datenschutz-Schutzziele [RB11] geben dieser Vielfalt Orientierung. Sie konkretisieren das Recht auf informationelle Selbstbestimmung ohne der Gestaltung des Datenschutzrechts im Einzelnen vorzugreifen. Im Einzelnen lauten die Datenschutzziele *Vertraulichkeit*, *Integrität*, *Verfügbarkeit*, *Transparenz*, *Unverkettbarkeit* und *Intervenierbarkeit*. Diese Ziele sind allerdings nicht konfliktfrei. Insbesondere Transparenz und Unverkettbarkeit stehen im Widerspruch zueinander [HJR15].

Das Ziel der *Transparenz* fußt auf dem Recht des Einzelnen, beurteilen zu können, wer wann was über ihn weiß.³ Das Auskunftsrecht als Teil der Transparenz ist für den Betroffenen das wichtigste Datenschutzrecht [Si11]. Es ist Voraussetzung zur Wahrnehmung

¹ Fraunhofer-Institut für Optronik, Systemtechnik und Bildauswertung IOSB, Fraunhoferstr. 1, 76131 Karlsruhe, christoph.bier@iosb.fraunhofer.de

² BVerfGE 65, 1 (43).

³ BVerfGE 65, 1 (43); BVerfGE 125, 260 (334).

der übrigen Betroffenenrechte auf Löschung, Sperrung und Berichtigung. Trotz seiner enormen Bedeutung für einen effektiven Datenschutz wird das Recht auf Auskunft durch die Praxis vernachlässigt. Auskünfte werden zwar erteilt, jedoch nur in Form von statischen Datenbankauszügen. An dieser Stelle setzt Personal-Data-Provenance an. Personal-Data-Provenance ist die dokumentierte Historie eines personenbezogenen Datums. Eine Provenance-Tracking-Infrastruktur verfolgt demnach den Lebenszyklus eines personenbezogenen Datums ausgehend von der Erhebung beim Betroffenen oder einem Dritten, über einzelne Verarbeitungsschritte bis hin zur Übermittlung. Alle Schritte werden mit dem Zweck der Erhebung und Verarbeitung des personenbezogenen Datums in Bezug gesetzt. Letztendlich soll der Betroffene die Möglichkeit bekommen, über eine Datenschutzauskunftsplattform jederzeit Einblick in den Umgang mit seinen personenbezogenen Daten zu nehmen [An15, BK16].

Unverkettbarkeit soll verhindern, dass es staatlichen Behörden und privaten Organisationen möglich ist, ein umfangreiches Persönlichkeitsprofil über jeden Einzelnen zu erstellen.⁴ Unverkettbarkeit fordert, dass personenbezogene Daten, Betroffene, Verarbeitungsprozesse und Nutzungsdomänen nicht miteinander verknüpft werden können. Aus datenschutzrechtlicher Sicht sind der Unverkettbarkeit insbesondere der Zweckbindungsgrundsatz⁵ sowie die informationelle und organisatorische Gewaltenteilung inhärent.

Transparenz erfordert eine ergänzende personenbezogene Sammlung von Protokolldaten und ihre Verknüpfung mit Bezug auf den Betroffenen. Jedes Mehr an Daten erhöht jedoch die Gefahr der Verkettbarkeit. Transparenz setzt die vollständige Verknüpfung von Daten und Verarbeitungsprozessen mit dem Betroffenen voraus, während Unverkettbarkeit des genauen Gegenteils bedarf. Das Ziel des in diesem Beitrag vorgestellten Ansatzes ist es, diesen Widerspruch sichtbar zu machen, um ihn, wenn möglich, aufzulösen. Entlang des Anwendungsfalls „Datenschutzauskunftssystem“ wird ein allgemeingültiges Modell für eine informationstheoretische Metrik für Unverkettbarkeit entworfen. Sie wird für ein Beispiel instantiiert und erläutert.

Verwandte Arbeiten. Ein Standardwerk zur Terminologie von Unverkettbarkeit in der Informatik ist bis heute die bis 2010 aktualisierte Veröffentlichung von Pfitzmann und Hansen [PH10]. Bohli und Pashalidis [BP11] formalisieren unterschiedliche Abstufungen der klassischen Unverkettbarkeit basierend auf der Unterscheidbarkeit von Systemrealisationen durch Angreifer mit unterschiedlichem, fest definierten Hintergrundwissen im Stil des IND-CPA-Modells zur Sicherheit kryptographischer Verfahren.

Die Idee, Anonymität informationstheoretisch zu beschreiben, wird bereits von Serjantov und Danezis ins Spiel gebracht [SD03]. Sie überführen das klassische „Anonymity Set“ auf ein nach den Wahrscheinlichkeiten der einzelnen Elemente der Menge gewichtetes Maß. Die Arbeiten von Steinbrecher und Köpsel [SK03] übertragen den informationstheoretischen Ansatz auf Unlinkability. Der Ansatz wird von Pashalidis [Pa08] von Äquivalenzrelationen auf alle zweistelligen Relationen verallgemeinert.

⁴ BVerfGE 65, 1 (42).

⁵ BVerfGE 65, 1 (43).

Gliederung. Der verbleibende Beitrag ist wie folgt gegliedert: Im Anschluss wird das Minimalbeispiel, anhand dessen die Berechnung der Metrik für Unverkettbarkeit verdeutlicht wird, geschildert. Anschließend werden im Abschnitt 1.1 die rechtlichen Anforderungen an die Ausgestaltung einer Unverkettbarkeitsmetrik erörtert. Das Systemmodell, die betrachteten Entitäten, das Angreifermodell und die informationstheoretische Metrik für Unverkettbarkeit werden in Abschnitt 2 eingeführt und im Abschnitt 2.3 für die unterschiedlichen Instanzierungen der Metrik konkret beschrieben. Ihre Implementierung wird im Abschnitt 2.4 erläutert. Im letzten Abschnitt wird der Nutzen der beschriebenen Metrik bewertet und ein Ausblick auf mögliche Erweiterungen gegeben.

Fortlaufendes Beispiel. Das fortlaufende Beispiel in diesem Beitrag nimmt ein Unternehmen namens AdBokis Buchclub GmbH, einen fiktiven Online-Händler für Bücher und Software, in den Fokus. Alice Fox ist Kundin dieses Händlers und möchte nach erfolgtem Einkauf ihr Auskunftsrecht wahrnehmen. Sie hat nach § 34 BDSG Anspruch auf Auskunft über die zu ihrer Person gespeicherten personenbezogenen Daten, deren Herkunft, Empfänger und den Zweck der Speicherung. Empfänger können der Betroffene, Dritte, Auftragsdatenverarbeiter und Stellen innerhalb der verantwortlichen Stelle sein. Der Stellenbegriff ist funktional und organisatorisch definiert [Si11].

In unserem Minimalbeispiel hat AdBokis zwei Kunden (Betroffene $b \in \mathcal{B}$): Alice Fox (b_1, s_1) und Peter Trollig (b_2, s_2). Die CloudyCloud GmbH (s_3) ist als Auftragsdatenverarbeiter für AdBokis tätig. Außerdem übermittelt AdBokis im Rahmen ihrer Geschäftsprozesse personenbezogene Daten an die PayPortal Inc. (s_4) und die Bonus Card GmbH (s_5). Intern spielen bei der Datenverarbeitung die Abteilungen Kundenbetreuung (s_6), Vertrieb, IT und Infrastruktur (s_{10}) und Recht (s_{11}) eine Rolle. In der Abteilung Vertrieb wird neben dem System für den Onlineverkauf (s_7) auch ein Archivserver (s_8) betrieben. Zudem gibt es Arbeitsplatzsysteme, die im Vertrieb normalerweise nicht für die Verarbeitung personenbezogener Daten vorgesehen sind. Exemplarisch ist deshalb im Minimalbeispiel der Workspace23 (s_9) enthalten. Alle diese Entitäten werden als Systeme $s \in \mathcal{S}$ bezeichnet. Die Verwendung des Begriffs *System* darf nicht verwirren. Ein System ist eine Ansammlung technischer (z.B. ein Cluster) oder organisatorischer (z.B. eine Abteilung) Entitäten, denen ein gemeinsames Wissen unterstellt wird.

Insgesamt verarbeitet AdBokis 30 personenbezogene Daten $d \in \mathcal{D}$ ihrer beiden Kunden in 17 Datenkategorien $\theta \in \Theta$. In Tabelle 1 sind exemplarisch die personenbezogenen Daten von Alice Fox aufgelistet.

Die Verarbeitung personenbezogener Daten findet entlang etablierter Verarbeitungsprozesse statt. Diese sind von der verantwortlichen Stelle gemäß § 4g Abs. 2 S. 1 BDSG i.V.m § 4e Satz 1 BDSG in einem internen Verzeichnisse zu dokumentieren. Teil dieser Dokumentation sind die verarbeiteten Datenkategorien, eine Beschreibung des Verfahrens sowie die möglichen Empfänger der Daten. Aus letzterer Information ergibt sich die Verknüpfung einzelner Verarbeitungsprozesse.

Auf eine Beschreibung des Datenschutzauskunftssystems wird aus Platzgründen an dieser Stelle verzichtet. Die wesentlichen Eigenschaften werden im Rahmen des fortlaufenden Beispiels erläutert.

d_x	θ_x	Datenkategorie	Inhalt
1	1	Vorname	Alice
2	2	Name	Fox
3	3	e-Mail	alice.fox@honigmail.de
		...	
16	14	Profilbild	[nicht darstellbar]
17	15	IP-Adresse	217.146.191.19
18	15	IP-Adresse	31.130.202.80
19	16	Empfehlung	Inges Braustubenführer
20	17	Rechnung	[nicht darstellbar]

Tab. 1: Auszug der Verarbeiteten Daten mit Personenbezug zu Alice

1.1 Aus der Unverkettbarkeit resultierende Anforderungen an ein Datenschutzauskunftssystem

Wie bereits im Abschnitt 1 erwähnt folgen aus der Unverkettbarkeit die Anforderungen der Zweckbindung und Zwecktrennung sowie die Organisatorische und technische Gewaltenteilung.

Zweckbindung und Zwecktrennung. Die Zweckbindung fordert, dass personenbezogene Daten nur zu dem Zweck verarbeitet und genutzt werden dürfen, zu dem sie auch erhoben wurden und der für sie dokumentiert wurde. Unter vielen findet sich diese Festlegung beispielsweise in § 28 Abs. 3 S. 7 BDSG. Die Zwecktrennung ist Ausfluss eines wesentlichen Aspekts der Unverkettbarkeit. Personenbezogene Daten, die zu unterschiedlichen Zwecken verarbeitet, insbesondere gespeichert, werden, dürfen nicht zusammengeführt werden. Mit technischem Bezug ist dies in der Anlage zu § 9 S. 1 BDSG Nr. 8 festgelegt.

Organisatorische und technische Gewaltenteilung. Die organisatorische und technische Gewaltenteilung im Rahmen der informationellen Gewaltenteilung folgt aus dem Gebot der Zweckbindung und Zwecktrennung. Ihr liegt das verwaltungsrechtliche Abschottungsprinzip zugrunde. Die informationelle Gewaltenteilung ist nicht im BDSG festgelegt, sondern ergibt sich aus der Rechtsprechung des Bundesverfassungsgerichts,⁶ in dessen Lichte das BDSG auszulegen ist. Während sich Zweckbindung und Zwecktrennung auf die personenbezogenen Daten selbst beziehen, ist die informationelle Gewaltenteilung eine Forderung, die direkt an die organisatorischen und technischen Einrichtungen gestellt wird. Die Zwecktrennung untersagt die Zusammenführung personenbezogener Daten. Die Gewaltenteilung verpflichtet dazu, dass Zugriffsrechte, Rollen sowie physische und logische Speicherorte nicht beliebig festgelegt werden. Sie sind entsprechend dem Zweck und nach dem Prinzip der Machtdistribution zu bestimmen.

Die organisatorische und technische Gewaltenteilung findet sich bisher noch kaum in expliziten Regelungsvorgaben des Bundesdatenschutzgesetzes. Sie ergibt sich jedoch indirekt aus anderen Anforderungen und Kriterien. Zunächst darf durch ein Datenschutzaus-

⁶ BVerfGE 65, 1 (69); BVerfG, NJW 1988, 959 (961).

kunftssystem keine Verknüpfung unterschiedlicher Daten möglich sein, soweit dies für die Auskunft nicht zwingend erforderlich ist. Außerdem dürfen die Bearbeitungswege und Weitergaben eines personenbezogenen Datums für eine Stelle, auch wenn sie an der Verarbeitung beteiligt ist, nicht durch das Datenschutzauskunftssystem nachvollziehbar gemacht werden. Daraus ergeben sich direkt die in Abschnitt 2.2 beschriebenen vier Unverkettbarkeitsrelationen.

2 Eine Metrik für Unverkettbarkeit

Unabdingbare Voraussetzung um über Unverkettbarkeit sprechen zu können, ist ein Modell des Ausschnitts der Realität zu definieren, in dem Unverkettbarkeit gemessen werden soll. Insbesondere sind die betrachteten Entitäten (\mathcal{E} - *items of interest*) und ihre Beziehungen zueinander (Verkettungsrelationen R) festzulegen. Außerdem ist ein Angreifer \mathcal{A} zu definieren, aus dessen Perspektive die Unverkettbarkeit bestimmt wird [PH10]. Eine Verkettungsrelation R ist eine Teilmenge des kartesischen Produkts von $n \geq 2$ Teilmengen $E_1, \dots, E_n \subseteq \mathcal{E}$ der Menge \mathcal{E} aller Entitäten: $R \subseteq E_1 \times \dots \times E_n$. (im Regelfall die Entitätsmengen der Entitätsklassen).

Unverkettbarkeit kann absolut oder relativ definiert werden. Absolut gesehen sind zwei oder mehrere Entitäten aus Sicht eines Angreifers dann unverkettbar, wenn der Angreifer nicht feststellen kann, ob die Entitäten innerhalb des definierten Modells in einem bestimmten Verhältnis zueinander stehen, oder nicht [PH10]. In einer Situation, in der ein gewisser Wissenszuwachs unabdingbar ist und in Kauf genommen wird, wie bei Auskunftssystemen, ist solch eine Metrik nicht hilfreich. Sie würde jederzeit trivial messen, dass die Unverkettbarkeit nicht gewahrt bleibt. Relative Unverkettbarkeit vergleicht die Unsicherheit des Angreifers \mathcal{A} bezüglich der wahren Verkettungsrelation R_τ nach Interaktion mit dem Gesamtsystem $\Sigma^{\mathcal{A}}$ mit der Unsicherheit, die bereits vor der Interaktion mit dem modellierten System bestand. Die Unsicherheit vor Interaktion ist vom Hintergrundwissen (A-priori-Wissen) des Angreifers abhängig. Die Interaktion mit $\Sigma^{\mathcal{A}}$ lässt den Angreifer die Beobachtungen I machen. Das kombinierte Wissen des Angreifers aus Hintergrundwissen und Beobachtungen wird auch als A-posteriori-Wissen bezeichnet. In der Literatur hat sich die informationstheoretische Bestimmung relativer Unverkettbarkeit etabliert (siehe Abschnitt 1). Im gesamten Beitrag werden Wahrscheinlichkeiten gemäß des Bayesschen Wahrscheinlichkeitsbegriffes als „Grad (vernünftiger) Glaubwürdigkeit/persönlicher Überzeugung“ (*degree of belief*) verwendet.

Sei X eine Zufallsvariable über der endlichen Menge der Kandidatenrelationen \mathcal{R} . Sowohl vor als auch nach Interaktion mit dem Gesamtsystem $\Sigma^{\mathcal{A}}$ weist der Angreifer \mathcal{A} allen Kandidatenrelationen $R \in \mathcal{R}$ einen Wahrscheinlichkeitswert $\mathbb{P}(X = R)$ zu. $\mathbb{P}(X = R)$ ist die angenommene Wahrscheinlichkeit, dass R die tatsächliche Relation R_τ zwischen den Entitäten aus E_1, \dots, E_n ist. Dann ergibt sich die Entropie des A-priori- bzw. A-posteriori-Wissens des Angreifers als:

$$H(X) = - \sum_{R \in \mathcal{R}} \mathbb{P}(X = R) \log_2 \mathbb{P}(X = R) \quad [\text{bit}]$$

Wobei $\mathbb{P}(X = R) \log_2 \mathbb{P}(X = R) = 0$ für $\mathbb{P}(X = R) = 0$ angenommen wird. Die Entropie misst die Informationsmenge, die \mathcal{A} noch braucht, um R_τ vollständig zu identifizieren.

Der Grad der Unverkettbarkeit ergibt sich als Verhältnis zwischen A-priori- und A-posteriori-Entropie (mit dem Beobachtungsereignis I):

$$\Delta(X, I) = \frac{H(X | I)}{H(X)}$$

Der Grad der Unverkettbarkeit beschreibt das Verhältnis zwischen der Situation nach und der Situation vor der Interaktion des Angreifers \mathcal{A} mit dem System $\Sigma^{\mathcal{A}}$ bezüglich des noch benötigten Wissens zur vollständigen Aufdeckung der Relation.

In bisherigen Arbeiten werden meist die A-priori-Situation und die maximale Unverkettbarkeit gleichgesetzt (Maximum-Entropie-Prior, $H(X) = H_{\max}(X)$). Die maximale Unverkettbarkeit ist $H_{\max}(X) = \log_2(|\mathcal{R}|)$. Ein Maximum-Entropie-Prior macht im diskutierten Szenario indes keinen Sinn. A priori sind damit Beobachtungen aus Datenverarbeitungsvorgängen ohne den Einsatz eines Datenschutzauskunftssystems vorausgesetzt. A posteriori werden die Beobachtungen aus den selben Datenverarbeitungsvorgänge unter Berücksichtigung des Einsatzes eines Datenschutzauskunftssystems ins Angreiferwissen mit aufgenommen. Statt eines Maximum-Entropie-Priors ist der Vergleichszustand (subjektiver Prior) schon ein Zustand mit partiellem Wissen.

Sei $H(X) \neq H_{\max}(X)$, dann ist bei Beobachtungen, die der A-priori-Annahme entgegengesetzt sind, hypothetisch ein $\Delta(X, I) > 1$ möglich. Allerdings ist als Maß der Unverkettbarkeit nicht die Grad der Unverkettbarkeit bezüglich eines bestimmten Angreifers von Interesse, sondern der niedrigste Grad über alle Angreifer. Trivial lässt sich immer ein Angreifer konstruieren, der kein Hintergrundwissen hat ($H(X) = H_{\max}(X)$) und durch seine Beobachtungen nichts dazulernen kann ($H(X | I) = H(X)$). Dessen Grad der Unverkettbarkeit ist immer $\Delta(X, I) = 1$. Somit ist der normierte globale Grad der Unverkettbarkeit $\|\Delta\| = \min_{\mathcal{A}} (\{\Delta(X_{\mathcal{A}}, I_{\mathcal{A}})\}) \in [0; 1]$.

2.1 Der Angreifer \mathcal{A}

Das Angreifermodell gibt die Leitlinien vor, an denen sich die Analyse der A-priori- und der A-posteriori-Situationen orientieren kann. Der Angreifer kann Teil der datenverarbeitenden Organisation (verantwortliche Stelle) sein oder außerhalb der Organisation zu finden sein. Im Folgenden werden exemplarisch nur zwei interne Angreifer betrachtet.

Der *Systemangreifer* \mathcal{A}^s verarbeitet möglicherweise selbst personenbezogene Daten. Er möchte aber Wissen über weitere Verarbeitungsvorgänge gewinnen. Vorstellbar ist beispielsweise eine Marketingabteilung, die wissen möchte, in welchem Maße und unter Preisgabe welcher Informationen ein Kunde bisher den Kundenservice angefragt hat. Der *zentrale Angreifer* \mathcal{A}^c entsteht erst durch das Datenschutzauskunftssystem. Die Datenschutzauskunft erfordert eine Aggregation der Personal-Data-Provenance vor der Weitergabe an den Betroffenen. Sie setzt voraus, dass ein Einstiegspunkt für den Abruf der

gesamten Provenance-Kette bekannt ist. Nur so kann die Vollständigkeit der Provenance gewährleistet werden. \mathcal{A}^s und \mathcal{A}^c werden als passive Angreifer angenommen. Sie halten die festgelegten Kommunikationsprotokolle des Datenschutzauskunftssystems vollständig ein. Eine Missachtung der Kommunikationsprotokolle kann von den Kommunikationspartnern festgestellt und organisatorisch verfolgt werden.

Das a-priori Hintergrundwissen der Angreifer umfasst das interne Verzeichnisse und allgemeine Unternehmensstatistiken. Das Verzeichnisse beinhaltet Informationen zu den Verarbeitungsprozessen, den beteiligten Systemen und den verwendeten Datenkategorien (siehe Abschnitt 1).

Annahme 1. *Dem Angreifer ist die Art und die Anzahl aller Systeme $s \in \mathcal{S}$ a-priori bekannt.*

Annahme 2. *Dem Angreifer ist die Anzahl der von der Datenverarbeitung betroffenen Kunden $|\mathcal{B}|$ a-priori bekannt. Dem Angreifer ist die Anzahl der verarbeiteten personenbezogenen Daten $|\mathcal{D}|$ a-priori bekannt.*

Personenbezogene Daten eines Betroffenen werden unabhängig davon erhoben und verarbeitet, ob personenbezogene Daten anderer Betroffener erhoben oder verarbeitet werden. Ob ein Kunde beispielsweise die Bezahlung mit Kreditkarte wählt hat keinen Einfluss darauf, ob dies ein anderer Kunde auch tut.

Annahme 3. *Die Zugehörigkeit eines personenbezogenen Datums zu einem Betroffenen (und umgekehrt) ist unabhängig von der Zugehörigkeit eines anderen personenbezogener Datums zu einem anderen Betroffenen (und umgekehrt).*

Ein Datum kann potentiell personenbezogenes Datum mehrerer Betroffener sein. Um die folgenden Überlegungen zu vereinfachen, wird dennoch angenommen, dass ein Datum nur einen Personenbezug zu einem Betroffenen haben kann.

Annahme 4. *Das Verhältnis von personenbezogenen Daten und Betroffenen ist eine $n:1$ -Beziehung.*

Die beiden letzten Annahmen 5 und 6, sind wichtige Annahmen zur Unabhängigkeit von Datenflüssen. Sie sind eine entscheidende Voraussetzung für die Berechenbarkeit der Unverkettbarkeitsmetriken. Beide Annahmen werden in der Realität nicht unbedingt vollständig eingehalten. Die durch sie induzierte Ungenauigkeit kann jedoch nur zu einem Unterschätzen des A-priori-Wissens des Angreifers führen. Das Delta zur A-posteriori-Situation wird dann größer. Die Gefährdung für den Datenschutz wird überschätzt. Deshalb sind die Annahmen vom Ergebnis her gedacht sinnvoller, als unbelegte Annahmen über die Abhängigkeit von Datenflüssen zu treffen, welche zu einem Unterschätzen des Datenschutzrisikos führen könnten.

Annahme 5. *Das a-priori Wissen zu Datenflüssen (Verarbeitungsprozessen) ist nur von der Kategorie der Daten, nicht von den Daten selbst abhängig.*

Annahme 6. *Die Flüsse zweier Daten sind unabhängig voneinander.*

2.2 Instanziierung der Unverkettbarkeit als Gegenspieler der Transparenz

Die relevanten Entitäten ergeben sich aus den Teilinformationen der Datenschutzauskunft. Es sind die Systeme $s \in \mathcal{S}$, die personenbezogenen Daten $d \in \mathcal{D}$ und die Betroffenen $b \in \mathcal{B}$. Gleiches gilt für die Verkettungsrelationen, die über diesen Entitäten definiert sind. Sie bilden das Interesse des Angreifers an den zu einem Betroffenen gespeicherten personenbezogenen Daten ($R^<$), an der Herkunft und den Empfängern personenbezogener Daten ($R^>$) und am zweckbestimmten Verarbeitungsort personenbezogener Daten (R^∇) ab. $R^=$ ergibt sich aus dem Gebot der Zwecktrennung. Wird die Zwecktrennung überwunden, kann ein Persönlichkeitsprofil des Betroffenen hergestellt werden, unabhängig davon, ob schon klar ist, wer er ist. Die vier genannten Relationen sind wie folgt definiert: (1) Die *Identifikationsrelation* $R^< \subseteq \mathcal{D} \times \mathcal{B}$ gibt an, ob das Datum $d \in \mathcal{D}$ einen Personenbezug auf den Betroffenen $b \in \mathcal{B}$ besitzt. (2) Die *Verknüpfungsrelation* $R^= \subseteq \mathcal{D} \times \mathcal{D}$ gibt an, ob zwei Daten $d_1, d_2 \in \mathcal{D}$ einen Personenbezug auf denselben (aber unbekanntem) Betroffenen besitzen. (3) Die *Speicher- und Verarbeitungsrelation* $R^\nabla \subseteq \mathcal{S} \times \mathcal{D}$ gibt für alle Systeme $s \in \mathcal{S}$ an, ob sie das Datum $d \in \mathcal{D}$ verarbeitet und/oder gespeichert haben. (4) Die *Datenflussrelation* $R^> \subseteq \mathcal{S} \times \mathcal{S} \times \mathcal{D}$ gibt für zwei Systeme $s_1, s_2 \in \mathcal{S}$ an, ob sie für ein bestimmtes personenbezogenes Datum $d \in \mathcal{D}$ in einer direkten Vorgänger-Nachfolger-Beziehung stehen.

$\Delta(X^<, I)$ und $\Delta(X^=, I)$ sind globale Metriken. Bei der Bestimmung des Grads der Unverkettbarkeit sind alle Betroffenen \mathcal{B} und alle personenbezogenen Daten \mathcal{D} mit einzubeziehen. Anders stellt sich die Situation bei $R^>$ und R^∇ dar. Der Grad der Unverkettbarkeit bezüglich dieser Mengen ist global und lokal bestimmbar. Lokal meint die Fokussierung auf bestimmte Systeme oder Betroffenen. Im Kontext der Datenschutzauskunft ist für einen Betroffenen nur relevant, wie sich die Unverkettbarkeit der Flüsse seiner personenbezogenen Daten entwickelt. Deshalb werden im Abschnitt 2.3.2 nur die Daten in der Teilmenge $\mathcal{D}_{\mathcal{B}} \subseteq \mathcal{D}$. Im Text wird dennoch im Sinne einer allgemeingültigen Darstellung von \mathcal{D} gesprochen. Gleichzeitig wird im Abschnitt 2.3.2 angenommen, dass dem Angreifer a priori bekannt ist, welche personenbezogenen Daten (aber nicht welcher Kategorie) zu welchem Betroffenen gehören. Die Unsicherheit über dieses Faktum wird bereits durch den Grad der Unverkettbarkeit von $R^<$ gemessen.

2.3 Modellierung des A-priori- und A-posteriori-Wissens der Angreifer

Um den Grad der Unverkettbarkeit bezüglich der vier genannten Relationen zu bestimmen, ist es erforderlich, das Hintergrundwissen der Angreifer und den Wissenszuwachs durch die Einführung der Datenschutzauskunft messbar zu machen. Das Hintergrundwissen der Angreifer geht in die A-priori-Wahrscheinlichkeiten $\mathbb{P}(X = R)$ mit ein. Der Wissenszuwachs der Angreifer wird durch das Beobachtungsereignis I und die daraus resultierenden A-posteriori-Wahrscheinlichkeiten $\mathbb{P}(X = R | I)$ erklärt. A-priori- und A-posteriori-Wahrscheinlichkeitsverteilungen aller vier Relationen werden in diesem Abschnitt erläutert und anhand des Minimalbeispiels aus Kapitel 1 bestimmt.

2.3.1 Bestimmung der Wahrscheinlichkeitsverteilungen von $X^<$ und $X^=$

Die Mächtigkeit der Menge der Kandidatenrelationen ist unter Berücksichtigung von Annahme 4 für die Identifikationsrelation durch $|\mathcal{R}^<| = |\mathcal{B}|^{|\mathcal{D}|}$ gegeben.

Beispiel. Die Anzahl der Kandidatenrelationen für $|\mathcal{D}| = 30$ und $|\mathcal{B}| = 2$ ist 2^{30} .

$\mathcal{R}^=$ ist eine Äquivalenzrelation. Bei Äquivalenzrelation ist die Anzahl möglicher Relationen durch die Bellsche Zahl $B_{|\mathcal{D}|}$ gegeben. Die Bellsche Zahl lässt sich mit Hilfe der Stirling-Zahl zweiter Art bestimmen. Da $\mathcal{R}^=$ auf $\mathcal{R}^<$ zurückzuführen ist, sind die tatsächlichen Kandidatenrelationen durch die Anzahl der Betroffenen $|\mathcal{B}|$ beschränkt. Nur k -Partitionen mit $k \leq |\mathcal{B}|$ sind möglich. Die Formel ist deshalb in korrigierter Form anzuwenden: $|\mathcal{R}^=| = \sum_{k=0}^{\min(|\mathcal{B}|, |\mathcal{D}|)} S_{|\mathcal{D}|, k}$

Beispiel. Die Anzahl der Kandidatenrelationen für $|\mathcal{D}| = 30$ und $|\mathcal{B}| = 2$ ist $|\mathcal{R}^=| = \sum_{k=0}^2 S(30, k) = 536870912$.

Für beide Relationen gilt, dass es, mit Ausnahme der Mächtigkeit der Mengen \mathcal{B} und \mathcal{D} , kein globales Hintergrundwissen gibt.

Für den Systemangreifer \mathcal{A}^s unterscheidet sich der Unverkettbarkeitsprior nicht vom Posterior. Durch das Datenschutzauskunftssystem werden auf den Systemen nur solche Teile der Provenance vorgehalten, die auf Ereignisse im jeweilige System zurückzuführen sind. Daraus folgt für den Systemangreifer $\Delta(X^<, I) = \Delta(X^=, I) = 1$. Der zentrale Angreifer verarbeitet selbst keine personenbezogenen Daten. Für ihn sind $H(X^<) = H_{\max}(X^<) = \log_2 |\mathcal{R}^<|$ und $H(X^=) = H_{\max}(X^=) = \log_2 |\mathcal{R}^=|$. A posteriori, also unter Einsatz des Datenschutzauskunftssystems, erhält der zentrale Angreifer weitere Informationen I . Bei jeder Erhebung eines personenbezogenen Datums wird ihm ein pseudonymer Identifikator für das erhobene Datum zusammen mit Informationen zum Betroffenen übermittelt. Daraus kann der Angreifer zwar nicht schließen, welches personenbezogene Datum oder welche Kategorie personenbezogener Daten erhoben wurde. Allerdings kann er bestimmen, wie viele personenbezogenen Daten für jeden einzelnen Betroffenen erhoben wurden. Kandidatenrelationen, die keine entsprechende Struktur aufweisen, kann er ausschließen. Seien die dem Angreifer bekannt gewordenen k -Partitionen für die Menge der Daten \mathcal{D} von der Größe l_1, l_2, \dots, l_k . Dann sind

$$|(R^= \in \mathcal{R}^= | \mathbb{P}(X^= = R^= | I) \neq 0)| = \binom{|\mathcal{D}|}{l_1} \binom{|\mathcal{D}| - l_1}{l_2} \dots \binom{|\mathcal{D}| - (l_1 + l_2 + \dots + l_{k-1})}{l_k}$$

und

$$|(R^< \in \mathcal{R}^< | \mathbb{P}(X^< = R^< | I) \neq 0)| = k! \cdot |(R^= \in \mathcal{R}^= | \mathbb{P}(X^= = R^= | I) \neq 0)|.$$

Unter Weiterbestehen der Gleichverteilungsannahme ergibt sich die A-posteriori-Entropie direkt aus der Mächtigkeit der obigen beiden Mengen.

Beispiel. Für Alice wurden 20 personenbezogene Daten erhoben, für Peter 10. Damit sinkt die Anzahl der möglichen Relationen $\mathcal{R}^=$ auf $|(R^= \in \mathcal{R}^= | \mathbb{P}(X^= = R^= | I) \neq 0)| = \binom{30}{20} =$

$\binom{30}{10} = 30045015$. Folglich ist der resultierende Grad der Unverkettbarkeit $\Delta(X^{\equiv}, I) = \frac{\log_2 30045015}{\log_2 2^{29}} \approx \frac{24,8406}{29} \approx 0,8566$.

Die Anzahl der verbleibenden Identifikationsrelation ist $|(R^{\leq} \in \mathcal{R}^{\leq} \mid \mathbb{P}(X^{\leq} = R^{\leq} \mid I) \neq 0)| = 21 \cdot 30045015 = 60090030$. Entsprechend ist der resultierende Grad der Unverkettbarkeit $\Delta(X^{\leq}, I) = \frac{\log_2 60090030}{\log_2 2^{30}} \approx \frac{25,8406}{30} \approx 0,8614$.

2.3.2 Bestimmung der Wahrscheinlichkeitsverteilungen von X^{\triangleright} und X^{∇}

Unter der Annahme, dass die Weitergaben unterschiedlicher personenbezogener Daten voneinander unabhängig sind (Annahme 6), kann die Wahrscheinlichkeit $\mathbb{P}(X^{\triangleright} = R^{\triangleright})$ für eine Kandidatenrelation $R^{\triangleright} \in \mathcal{R}^{\triangleright}$ aus den Wahrscheinlichkeiten für die Teilrelationen je Datum $\mathbb{P}(X_d^{\triangleright} = R_d^{\triangleright})$ mit $R_d^{\triangleright} \subseteq \mathcal{S} \times \mathcal{S}$ berechnet werden. Es gilt $\mathbb{P}(X^{\triangleright} = R^{\triangleright}) = \prod_{d \in \mathcal{D}} \mathbb{P}(X_d^{\triangleright} = R_d^{\triangleright})$.

Das Wissen eines Angreifers wird zunächst dadurch charakterisiert, inwiefern ihm die Kategorie des personenbezogenen Datums bekannt ist. Die Kategorie des Datums bestimmt dessen Herkunft und Verarbeitung gemäß Verfahrensverzeichnis. Jedem Datum ist seine Datenkategorie über die Funktion $\vartheta : \mathcal{D} \rightarrow \Theta$ zugewiesen.

$$\mathbb{P}(X^{\triangleright} = R^{\triangleright}) = \sum_{\theta \in \Theta} \mathbb{P}(X_d^{\triangleright} = R_d^{\triangleright} \mid X_\theta = \theta) \mathbb{P}(X_\theta = \theta)$$

Beispiel. Die Zuordnung zwischen Daten und Datenkategorien ist einem Systemangreifer für diejenigen Daten bekannt, die er selbst verarbeitet. So ist \mathcal{A}_{s_3} bekannt, dass $\vartheta(d_{16}) = \theta_{14}$ ist. Für alle anderen Daten sowie grds. für den zentralen Angreifer muss entsprechend der Maximum-Likelihood-Methode die Gleichverteilung angenommen werden.

Das Wissen eines Angreifers wird außerdem dadurch charakterisiert, inwiefern ihm die Herkunft der personenbezogenen Daten bekannt ist. Die Herkunft der personenbezogenen Daten ist von der Kategorie der Daten abhängig. Die Wahrscheinlichkeiten der Zufallsvariable X_σ für bestimmte Startsysteme abhängig von der Datenkategorie lautet $\mathbb{P}(X_\sigma = s \mid X_\theta = \theta)$. Und es ergibt sich:

$$\mathbb{P}(X_d^{\triangleright} = R_d^{\triangleright} \mid X_\theta = \theta) = \sum_{s \in \mathcal{S}} \mathbb{P}(X_d^{\triangleright} = R_d^{\triangleright} \mid X_\sigma = s, X_\theta = \theta) \mathbb{P}(X_\sigma = s \mid X_\theta = \theta)$$

Beispiel. Alle personenbezogenen Daten bis auf jene der Kategorien Rechnung und Empfehlung werden ausschließlich beim Betroffenen selbst erhoben. Für diese ist das Herkunftssystem s_1 bekannt. Eine Rechnung wird immer im Onlineshopsystem des Vertriebs erzeugt. Das Herkunftssystem für die Rechnung ist damit aus Sicht der Angreifer eindeutig s_7 . Über eine Empfehlung ist hingegen nur bekannt, dass sie von einem Kunden kommen muss. Der Startvektor für die Empfehlung ist mit der Wahrscheinlichkeit $\frac{1}{2} s_1$ und mit der selben Wahrscheinlichkeit s_2 .

Anmerkung Jede zweistellige⁷ Relation R über endlichen Mengen kann als binäre bzw. boolesche Matrix $R^{\square} = (r_{ij})$, $r_{ij} \in \{0, 1\}$ dargestellt werden. Die Einträge der Matrix r_{ij} stehen für die Realisationen der wie folgt definierten Zufallsvariablen X_{ij} :

$$X_{ij} : \mathcal{R} \rightarrow \{0, 1\}$$

$$R \mapsto \begin{cases} 1 & (e_i, e_j) \in R \\ 0 & \text{sonst.} \end{cases}$$

Der Eintrag bzw. das Ereignis 1 bedeutet, dass die die durch den Index gegebenen Elemente in Relation zueinander stehen, der Eintrag 0, dass keine Beziehung vorliegt. Im Folgenden wird deshalb zur Vereinfachung nicht zwischen der Relation R und ihrer Matrixdarstellung R^{\square} differenziert.

$$r_{ij} = 1 \Leftrightarrow X_{ij} = 1 \Leftrightarrow (e_i, e_j) \in R$$

Daraus abgeleitet wird folgende Kurzschreibweise verwendet:

$$R_{r_{i_1, j_1}, r_{i_2, j_2}, \dots, r_{i_k, j_k}} : \Leftrightarrow R = \{e_{i_1}, e_{j_1}\} \cup \{e_{i_2}, e_{j_2}\} \cup \dots \cup \{e_{i_k}, e_{j_k}\}$$

Bezüglich der konkreten Datenflüsse stützten sich die Angreifer auf die Angaben des Verfahrensverzeichnis. Dieses hinterlegt für alle Daten die vorgesehenen Verarbeitungsprozesse. Das Wissen der Angreifer wird als Matrix der bedingten Flusswahrscheinlichkeiten W modelliert. Der Eintrag w_{ij} gibt die Wahrscheinlichkeit an, mit der ein Fluss von s_i nach s_j , angenommen wird, unter der Bedingung, dass das Datum bereits in s_i verarbeitet wurde:

$$W_{\theta, s} : \{1, \dots, m\} \times \{1, \dots, m\} \rightarrow [0, 1]$$

$$(i, j) \mapsto w_{ij} = \mathbb{P}(X_{dij}^{\triangleright} = 1 \mid X_{dii}^{\triangleright} = 1, X_{\sigma} = s, X_{\theta} = \theta)$$

mit $m = |\mathcal{S}|$. Implizit ist $w_{ii} = 1$. Der reflexive Fluss, gleichbedeutend mit der Speicherung und Verarbeitung im System ($\forall d, i : \mathbb{P}(r_{dii}^{\triangleright}) = \mathbb{P}(r_{dii}^{\nabla})$), ist vollständig durch die eingehenden Flüsse erklärt:

$$\mathbb{P}(X_{dii}^{\triangleright} = 1 \mid \exists j \in \{1, \dots, i-1, i+1, \dots, m\} : X_{dji}^{\triangleright} = 1, X_{\sigma} = s, X_{\theta} = \theta) = 1$$

$$\mathbb{P}(X_{dii}^{\triangleright} = 1 \mid \forall j \in \{1, \dots, i-1, i+1, \dots, m\} : X_{dji}^{\triangleright} = 0, X_{\sigma} = s, X_{\theta} = \theta) = 0$$

Für alle Datenflussrelationen gilt $\mathbb{P}(X_{dij}^{\triangleright} = 0 \mid X_{dii}^{\triangleright} = 0, X_{\sigma} = s, X_{\theta} = \theta) = 1$ und $\mathbb{P}(X_{dij}^{\triangleright} = 1 \mid X_{dii}^{\triangleright} = 0, X_{\sigma} = s, X_{\theta} = \theta) = 0$. Es kann keine ausgehenden Flüsse geben, falls es keinen eingehenden Fluss gibt. Damit ist der Wahrscheinlichkeitsbaum für die Datenflussrelation vollständig erklärt.

Die in den bedingten Flusswahrscheinlichkeiten zum Ausdruck kommenden Pfade ergeben sich aus dem Verfahrensverzeichnis. Das Erfahrungswissen des Angreifers lässt sich in zwei zentralen Parametern ausdrücken. Zunächst hängt die Wahrscheinlichkeit von Flüssen in linearen Verfahren maßgeblich von der *Fortschrittsquote* $\varpi \in (0, 5; 1]$ eines

⁷ Gilt grundsätzlich auch für mehrstellige Relationen.

Prozesses ab. Dieser Parameter wird überall dort in der Flussmatrix eingesetzt, wo ein Fluss einem Prozessschritt entspricht. Unvorhergesehene Abweichungen vom Verfahren werden durch eine Fehlerwahrscheinlichkeit $\varepsilon \in [0; 0,5)$ beschrieben. Mit dieser Fehlerwahrscheinlichkeit finden Flüsse zu und zwischen Systemen außerhalb des vorgesehenen Prozessablaufs statt.

Beispiel. Die AdBokis Buchclub GmbH hat folgende Verfahren etabliert: Registrierung, Bestellung, Zahlungsabwicklung, Kundendatenarchivierung, Missbrauchsbekämpfung und Kundenservice. Die Fortschrittsquote wird mit 90% ($\varpi = 0,9$) und die Fehlerwahrscheinlichkeit mit $\varepsilon = 0,02$ angenommen. Auf die Darstellung der einzelnen Prozesse muss aus Platzgründen verzichtet werden.

Auf dieser Grundlage kann die A-priori-Wahrscheinlichkeit für die einzelnen Kandidatenrelationen iterativ berechnet werden. Die Komplexität der vollständigen Berechnung der Wahrscheinlichkeiten aller möglichen Kandidatenrelationen ist allerdings in $O(2^{|\mathcal{S}| \cdot |\mathcal{S}|})$. Selbst bei wenigen Systemen ist somit die Berechenbarkeit der Wahrscheinlichkeitsverteilung nicht mehr gegeben. Deshalb ist nur eine heuristische Lösung entsprechend dem in Abschnitt 2.4 beschriebenen Verfahren möglich.

Beispiel. Für obige Flussmatrix ergibt sich die Wahrscheinlichkeit, dass das Datum nur im Herkunftssystem s_1 verarbeitet wird, sofern es von der Datenkategorie θ_{14} ist als

$$\mathbb{P}(X_d^\triangleright = R_{d,r_{d11}^\triangleright=1}^\triangleright \mid X_\sigma = s_1, X_\theta = \theta_{14}) = 1 \cdot 0,1 \cdot 0,98^9 \approx 0,0834$$

Würde man kombinatorisch aus den Wahrscheinlichkeiten $\mathbb{P}(X_d^\triangleright = R_d^\triangleright)$ die Gesamtwahrscheinlichkeit $\mathbb{P}(X^\triangleright = R^\triangleright) = \prod_{d \in \mathcal{D}} \mathbb{P}(X_d^\triangleright = R_d^\triangleright)$ berechnen, hätte dies eine Komplexität in $O(|\mathcal{D}|^{|\mathcal{S}| \cdot |\mathcal{S}|})$. Erfreulicherweise gilt für unabhängige Teilsysteme (Teilrelationen), dass die Entropie eine additive Größe ist ($H(X^\triangleright) = H(X_{d_1}^\triangleright) + \dots + H(X_{d_n}^\triangleright)$ mit $n = |\mathcal{D}|$). Somit lässt sich die Gesamtwahrscheinlichkeit aus den approximierten Teilwahrscheinlichkeiten bestimmen.

Die Wahrscheinlichkeitsverteilung für X^∇ lässt sich auf Grundlage der Wahrscheinlichkeitsverteilung von X^\triangleright ermitteln:

$$\mathbb{P}(X^\nabla = R^\nabla) = \sum_{R^\triangleright \in \mathcal{R}^\triangleright \mid R^\triangleright \equiv_\nabla R^\nabla} \mathbb{P}(X^\triangleright = R^\triangleright)$$

mit $R^\triangleright \equiv_\nabla R^\nabla \Leftrightarrow$

$$\forall d \in \mathcal{D}, s \in \mathcal{S} : ((d, s, s) \in R^\triangleright \wedge (d, s) \in R^\nabla) \vee ((d, s, s) \notin R^\triangleright \wedge (d, s) \notin R^\nabla)$$

Zu diesem allgemeinen Hintergrundwissen kommen noch die jeweiligen Beobachtungen I der Angreifer hinzu. Ein Systemangreifer kann die Datenflüsse durch sein System überwachen. Die Likelihood $\mathbb{P}(I \mid R^\triangleright)$ ist für solche Beobachtungen sicher 1 oder 0. Die A-posteriori-Wahrscheinlichkeit beträgt

$$\mathbb{P}(R^\triangleright \mid I) = \frac{\mathbb{P}(I \mid R^\triangleright) \mathbb{P}(R^\triangleright)}{\mathbb{P}(I)} = \frac{\mathbb{P}(I \mid R^\triangleright) \mathbb{P}(R^\triangleright)}{\sum_{R^{\triangleright'} \in \mathcal{R}^\triangleright} \mathbb{P}(I \mid R^{\triangleright'}) \mathbb{P}(R^{\triangleright'})}$$

und damit entweder 0 oder $\frac{\mathbb{P}(R^{\triangleright})}{\sum_{R^{\triangleright} \in \mathcal{R}^{\triangleright}} \mathbb{P}(I, R^{\triangleright} I)}$. Es findet also eine Normierung auf die Summe der Wahrscheinlichkeiten der Relationen, die die Beobachtung des Angreifers zulassen, statt.

Der zentrale Angreifer kann ohne das Datenschutzauskunftssystem keine Beobachtungen machen, sondern muss sich vollständig auf das Hintergrundwissen auf Grundlage des Verfahrensverzeichnis verlassen. Er ist jedoch der einzige Angreifer, der mit Hilfe des Datenschutzauskunftssystems weitere Beobachtungen I' machen kann. Für \mathcal{A}^S ist $\Delta(X^{\nabla}, I') = \Delta(X^{\triangleright}, I') = 1$. Bei der Registrierung neu erhobener personenbezogener Daten im zentralen Verzeichnis lernt der zentrale Angreifer die Quelle der personenbezogenen Daten und den Ort der ersten Verarbeitung im Unternehmen kennen (Likelihood von 1 oder 0). Als Ergebnis lassen sich jeweils die A-posteriori-Wahrscheinlichkeiten und der abgeleitete Grad der Unverkettbarkeit für X^{\triangleright} und X^{∇} nach dem im Abschnitt 2.4 beschriebenen Verfahren bestimmen.

2.4 Implementierung

Wie bereits im vorherigen Abschnitt erwähnt, ist die Wahrscheinlichkeitsverteilung für die Relation R_d^{\triangleright} auch schon bei wenigen Systemen, Daten und Datentypen nicht mit akzeptablem Aufwand an Zeit und Speicher vollständig berechenbar. Allerdings ist eine approximative Lösung möglich. Sind die Wahrscheinlichkeitsmatrizen nur spärlich mit Fortschrittswahrscheinlichkeiten belegt, ballt sich die Wahrscheinlichkeitsmasse bei denjenigen Kandidatenrelationen, die einen Fluss entlang des Verarbeitungsprozesses vorsehen. Kandidatenrelationen, die kaum Flüsse im Verarbeitungsprozess vorsehen, bilden den „Long Tail“ der Verteilung. Ihr Gewicht bei der Berechnung der Entropie ist gering. Diesen Umstand kann man sich bei der Berechnung des *belief* aus den Wahrscheinlichkeitsmatrizen zunutze machen, indem man systematisch zuerst die wahrscheinlicheren Kandidaten in die Berechnung aufnimmt und den „Long Tail“ nur bis zu einem gegebenen *Schwellwert* erschließt. Die Wahrscheinlichkeiten ergeben sich aus dem Wahrscheinlichkeitsbaum. Durch Tiefensuche in diesem Baum kann die Entropie, ausgehend vom Startsystem, approximativ erschlossen werden.

Beispiel. *Tabelle 2 enthält die berechneten Ergebnisse für unterschiedliche Schwellwerte. Dem Betroffenen könnte ein Mindestgrad an Unverkettbarkeit von 0,94 bzw. 0,90 garantiert werden.*

3 Zusammenfassung und Ausblick

In diesem Beitrag wurde eine informationstheoretische Metrik für Unverkettbarkeit entworfen und für ein Datenschutzauskunftssystem instanziiert. Die Modellierung des Angreiferwissens wurde anhand eines Beispiels eingeführt. Die Metrik beruht auf definierten rechtlichen Anforderungen und passt sich in die Datenschutz-Schutzziele ein. Sie erlaubt, den Trade-off zwischen Transparenz und Unverkettbarkeit deutlich zu machen.

Schwellwert	$H(X^{\triangleright})$	$H(X^{\triangleright}, I')$	$\Delta(X^{\triangleright}, I')$	$H(X^{\nabla})$	$H(X^{\nabla}, I')$	$\Delta(X^{\nabla}, I')$
10^{-1}	125,8201	114,9513	0,8927	96,2104	85,8917	0,9136
10^{-2}	139,5569	128,8865	0,9235	103,6033	92,4655	0,8925
10^{-3}	149,4219	139,3821	0,9328	107,1627	96,1635	0,8974
10^{-4}	156,1207	147,2097	0,9429	109,1815	98,5819	0,9029
10^{-5}	157,3259	148,4618	0,9436	109,4951	98,8992	0,9032

Tab. 2: Approximierte Werte für die Entropie und den Grad der Unverkettbarkeit der Relationen R^{\triangleright} und R^{∇} für unterschiedliche Schwellwerte

Im Rahmen einer Nutzerevaluation zu unserem Datenschutzauskunftssystem [BK16] haben wir auch die Einstellung von Betroffenen zur Metrik abgefragt, die in Form von Zustandsbalken in das Frontend des Systems eingebunden war. Von unseren 31 Teilnehmern hatten 21 das Konzept verstanden, 7 wahren sich aufgrund der gegebenen Kurzbeschreibung nicht sicher. Von diesen 21 hielten nur 6 die Metrik nicht für Hilfreich, für 9 war es eine nützliche Entscheidungsgrundlage für ihre Opt-out-Möglichkeit. Insgesamt ein zufriedenstellendes Ergebnis für eine bis dato unbekannte Entscheidungshilfe.

Literaturverzeichnis

- [An15] Angulo, Julio; Fischer-Hübner, Simone; Pulls, Tobias; Wästlund, Erik: Usable Transparency with the Data Track: A Tool for Visualizing Data Disclosures. In: Proc. of the 33rd ACM Conference on Human Factors in Computing Systems. S. 1803–1808, 2015.
- [BK16] Bier, Christoph; Kühne, Kay: PrivacyInsight: The Next Generation Privacy Dashboard. In: Proc. of the Annual Privacy Forum (APF 2016). 2016.
- [BP11] Bohli, Jens-Matthias; Pashalidis, Andreas: Relations among privacy notions. ACM Transactions on Information and System Security, 14(1):1–24, 2011.
- [HJR15] Hansen, Marit; Jensen, Meiko; Rost, Martin: Protection Goals for Privacy Engineering. In: Proc. of the 1st Int. Workshop on Privacy Engineering. S. 159–166, 2015.
- [Pa08] Pashalidis, Andreas: Measuring the Effectiveness and the Fairness of Relation Hiding Systems. In: Proc. of the Asia-Pacific Services Computing Conf. S. 1387–1394, 2008.
- [PH10] A terminology for talking about privacy by data minimization: Anonymity, Unlinkability, Undetectability, Unobservability, Pseudonymity, and Identity Management. http://dud.inf.tu-dresden.de/literatur/Anon_Terminology_v0.34.pdf, v0.34.
- [RB11] Rost, Martin; Bock, Kirsten: Privacy By Design und die Neuen Schutzziele. Datenschutz und Datensicherheit, 35(1):30–35, 2011.
- [SD03] Serjantov, Andrei; Danezis, George: Towards an Information Theoretic Metric for Anonymity. In: Proc. 2nd Int. Conf. on Privacy Enhancing Technologies. S. 41–53, 2003.
- [Si11] Simitis, Spiros, Hrsg. BDSG. Nomos, Baden-Baden, 7. Auflage, 2011.
- [SK03] Steinbrecher, Sandra; Köpsell, Stefan: Modelling Unlinkability. In: Proc. of the 3rd Int. Workshop on Privacy Enhancing Technologies (PET 2003). Springer, S. 32–47, 2003.