

Epistemic Privacy

Alexandre Evfimievski Ronald Fagin David Woodruff

IBM Almaden Research Center
650 Harry Rd., San Jose, CA, USA

{evfimi, fagin, dpwoodru}@us.ibm.com

ABSTRACT

We present a novel definition of privacy in the framework of offline (retroactive) database query auditing. Given information about the database, a description of sensitive data, and assumptions about users' prior knowledge, our goal is to determine if answering a past user's query could have led to a privacy breach. According to our definition, an audited property A is private, given the disclosure of property B , if no user can gain confidence in A by learning B , subject to prior knowledge constraints. Privacy is not violated if the disclosure of B causes a loss of confidence in A . The new notion of privacy is formalized using the well-known semantics for reasoning about knowledge, where logical properties correspond to sets of possible worlds (databases) that satisfy these properties. Database users are modelled as either possibilistic agents whose knowledge is a set of possible worlds, or as probabilistic agents whose knowledge is a probability distribution on possible worlds.

We analyze the new privacy notion, show its relationship with the conventional approach, and derive criteria that allow the auditor to test privacy efficiently in some important cases. In particular, we prove characterization theorems for the possibilistic case, and study in depth the probabilistic case under the assumption that all database records are considered a-priori independent by the user, as well as under more relaxed (or absent) prior-knowledge assumptions. In the probabilistic case we show that for certain families of distributions there is no efficient algorithm to test whether an audited property A is private given the disclosure of a property B , assuming $P \neq NP$. Nevertheless, for many interesting families, such as the family of product distributions, we obtain algorithms that are efficient both in theory and in practice.

Categories and Subject Descriptors: H.2.7 [Database Management]: Database Administration; F.2.1 [Analysis of Algorithms and Problem Complexity]: Numerical Algorithms and Problems

General Terms: Algorithms, Security, Theory

Keywords: privacy, disclosure, auditing, query logs, reasoning about knowledge, supermodularity, Positivstellensatz

1. INTRODUCTION

Today, privacy protection has become a popular and even fashionable area of database research. This situation is, of course, quite natural, given the importance of privacy in our social life and the risks we face in the digital world. These risks were highlighted by numerous recent reports of personal data theft and misappropriation, prompting many countries to enact data protection laws. However, the current state of scientific knowledge still does not allow the implementation of a comprehensive privacy solution that guarantees provable protection. In fact, the notion of privacy itself

has many definitions and interpretations, some focused on theoretical soundness, others on practical usefulness. This paper attempts to reduce the gap between these two aspects by exploring more flexible yet sound definitions.

One typical privacy enforcement problem, called *query auditing*, is to determine if answering a user's database query could lead to a privacy breach. To state the problem more accurately, we assume that the auditor is given:

- The database at the time of the user's query, or some partial knowledge about that database;
- A description of information considered sensitive, often called the *privacy policy* or the *audit query*;
- Assumptions about the user's prior knowledge of the database, of the audit query / privacy policy, and of the auditor's privacy enforcement strategy if it exists;
- The user's query, or a range of queries.

The auditor wants to check whether answering a given query could augment the user's knowledge about some sensitive data, thereby violating the privacy of that data. This problem has two extensions: *proactive* privacy enforcement (also called *online auditing* [18]), and *retroactive* or *offline* auditing.

In the proactive (online) privacy enforcement scenario, users issue a stream of queries, and the database system decides whether to answer or to deny each query. The denial, when it occurs, is also an "answer" to some (implicit) query that depends on the auditor's privacy enforcement strategy, and therefore it may disclose sensitive data. The strategy has to be chosen in advance, before the user's queries become available. A strategy that protects privacy for a specified range of queries represents a solution to this auditing problem. An in-depth discussion of online auditing can be found in [18, 23] and papers referenced therein.

In the retroactive (offline) scenario, the users issue their queries and receive the answers; later, an auditor checks if a privacy violation might have occurred. The audit results are not made available to the users, so the auditor's behavior no longer factors into the disclosure of data, and this considerably simplifies the problem. This also allows for more flexibility in defining sensitive information: while in the proactive case the privacy policy is typically fixed and open to the users, in the retroactive case the audit query itself may be sensitive, e. g. based on an actual or suspected privacy breach [1, 22]. Retroactive auditing is the application that motivates this paper, although our framework turns out to be fairly general.

To further illustrate the above, suppose Alice asks Bob for his HIV status. Assume that Bob never lies and considers "HIV-positive" to be sensitive information, while "HIV-negative" is for him OK to disclose. Bob is HIV-negative at the moment; can he adopt the proactive strategy of answering "I am HIV-negative" as long as it is true? Unfortunately, this is not a safe strategy, because if he does become HIV-positive in the future, he will have to deny

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

PODS'08, June 9–12, 2008, Vancouver, BC, Canada.

Copyright 2008 ACM 978-1-60558-108-8/08/06 ...\$5.00.

further inquiries, and Alice will infer that he contracted HIV. The safest bet for Bob is to always refuse an answer.¹

For the retroactive scenario, suppose that Bob contracted HIV in 2006. Alice, Cindy and Mallory legitimately gained access to Bob’s health records and learned his HIV status, but Alice and Cindy did it in 2005 and Mallory did in 2007. Bob discovers that his disease is known to the drug advertisers, and he initiates an audit, specifying “HIV-positive” as the audit query. The audit will place the suspicion on Mallory, but not on Alice and Cindy.

In legal practice, retroactive law enforcement has shown to be better suited to the complex needs of our society, although proactive measures are used too, especially in simple or critical situations. For example, a valuable item can be protected from theft by lock and key (a proactive measure) or by the fear of being caught and jailed (a retroactive measure). If it is simple to fence off the item and distribute the keys to all authorized users, or if the item has extraordinary value, then proactive defense is the best option, but in less clear-cut cases this would be too cumbersome or intrusive. After all, even an authorized user might steal or lose the item, and even a stranger sometimes should be able to gain access to it, e. g. in an emergency. Healthcare [2] is one area where the complexity of data management is just too high to hope for a fully proactive solution to privacy. The importance of offline disclosure auditing in healthcare has been recognized by the U.S. President’s Information Technology Advisory Committee [26], which recommended that healthcare information systems have the capability to audit who has accessed patient records. We believe in coexistence and importance of both auditing approaches.

1.1 Privacy Definitions in Query Auditing

The art of encryption and cryptanalysis goes back to antiquity, but the scientific maturity of privacy theory was made possible only in modern times by mathematical modeling of the eavesdropper’s knowledge. One of the first such models was proposed in 1949 by Claude Shannon [29], who introduced the notion of *perfect secrecy*. Shannon suggested to represent the adversarial knowledge by a probability distribution over possible private data values: prior distribution before the cryptogram is revealed, and posterior distribution after the adversary sees the cryptogram (but not the key). Perfect secrecy corresponds to the situation where the posterior distribution is identical to the prior, for every possible cryptogram. This general idea has been later adapted and extended to many privacy frameworks and problems, including query auditing.

Denote by Ω the set of all possible databases, and by A and B two properties of these databases; each database $\omega \in \Omega$ either has or does not have each property. Assume that the actual database satisfies both A and B . Suppose that property A is sensitive, and property B is what user Alice has learned by receiving the answer to her query. Was the privacy of A violated by the disclosure of B ? This depends on what Alice knew before learning B ; for example, if she knew “ $B \Rightarrow A$ ” (but did not know A), then B of course revealed A to her. Miklau and Suciu [21] applied Shannon’s model to this problem and declared A to be private given B if and only if

$$P[A|B] = P[A] \quad (1)$$

for all probability distributions P over Ω that might describe Alice’s prior knowledge about the database. Unfortunately, if no constraints are placed on P , no pair (A, B) of non-trivial properties will satisfy this privacy definition. Miklau and Suciu considered a quite limiting, yet popular, constraint: that Alice treats all database records $r \in \omega$ independently, i. e. P is a product distribution:

$$P(\omega) = \prod_{r \in \omega} P[r] \times \prod_{r \notin \omega} (1 - P[r])$$

¹If Alice pays Bob for answers, he can balance privacy and profit by tossing a coin and answering “I am HIV-negative” only if the coin falls heads.

Under this constraint, they prove that property A is private given the disclosure of B if and only if they share no *critical records* (Theorem 3.5 in [21]). A database record, real or imaginary, is called “critical” for A (for B) if its presence or absence in some database may decide the truth value for A (for B). One can see that, even with prior knowledge restricted to product distributions, very few practical queries would get privacy clearance: usually we can find an imaginary record and a pair of imaginary databases, ω_A for A and ω_B for B , where inserting the record into ω_A (into ω_B) flips the truth value of A (of B). Perfect secrecy appears too demanding to be practical.

A number of recent papers studied ways to relax condition (1) and make it approximate. They follow the same principle: for certain pairs of numerical bounds (ρ_1, ρ_2) , $\rho_1 < \rho_2$, require that

$$P[A] \leq \rho_1 \Rightarrow P[A|B] \leq \rho_2$$

where P is a prior knowledge distribution. This idea is behind the definition of ρ_1 -to- ρ_2 privacy breach in [12]; Kenthapadi *et al.* [18] use a slightly different version as part of their definition:

$$1 - \lambda \leq P[A|B] / P[A] \leq 1 / (1 - \lambda)$$

The Sub-Linear Queries (SuLQ) framework developed in [5, 10, 11] has a more sophisticated version with nice theoretical characteristics:

$$\Pr \left[\log \frac{P[A|B]}{1 - P[A|B]} - \log \frac{P[A]}{1 - P[A]} > \varepsilon \right] \leq \delta \quad (2)$$

While there is no space here to thoroughly review these frameworks, conceptually they all require that no user can gain much confidence in the audited property A by learning the disclosed property B , subject to prior knowledge constraints.

Perhaps surprisingly, however, all papers known to us, in their proofs if not in their definitions, do not make any distinction between *gaining* and *losing* the confidence in A upon learning B . For example, the SuLQ results remain in force if the privacy definition of [5] is changed by placing the absolute value sign “[...]” over the difference in (2). In some papers [11] the “[...]” appears in the definition explicitly.

It turns out that taking advantage of the gain-vs.-loss distinction yields a remarkable increase in the flexibility of query auditing. To bring it into focus, we shall put aside the approximate privacy relaxations and replace Eq. (1) with inequality

$$P[A|B] \leq P[A] \quad (3)$$

That is, we call property A *private* given the disclosure of property B when (3) holds for all distributions P that are admissible as a user’s prior knowledge. One might call this “semiperfect secrecy,” for it has the same sort of “absolute” form as perfect secrecy. This and related notions are the subject of this paper.

Let us illustrate its flexibility with a simple example of Alice (a user) and Bob (a patient). The hospital’s database ω has two records: $r_1 =$ “Bob is HIV-positive” and $r_2 =$ “Bob had blood transfusions.” The sensitive property A is the presence of r_1 , i. e. that Bob is HIV-positive. The property B that Alice queries and learns is “ $r_1 \in \omega$ implies $r_2 \in \omega$,” i. e. that “if Bob is HIV-positive, then he had blood transfusions.” *We make no constraints on Alice’s prior knowledge distribution*, other than a nonzero probability of the actual database. Could the disclosure of B violate the privacy of A ? Look at the following table of possible worlds:

	$r_2 \in \omega$	$r_2 \notin \omega$
$r_1 \in \omega$	A is true	A is true ★
$r_1 \notin \omega$	A is false	A is false

For Alice, learning B has the effect of ruling out the cell marked with a \star , while leaving the other cells untouched. Whatever the cells' prior probabilities are, the odds of A can only go down: $P[A|B] \leq P[A]$. Thus, A is private with respect to B , even though A and B share a critical record r_1 , and regardless of any possible dependence among the records.²

1.2 Summary of Results

This paper studies a notion of database privacy that makes it illegal for users to gain confidence about sensitive facts, yet allows arbitrary confidence loss. We begin in Sections 2 and 3 by introducing two novel privacy frameworks that implement the above concept for two different knowledge representations: possibilistic and probabilistic. We outline some properties of our privacy definitions that are relevant to the problem of testing privacy, and give necessary and sufficient conditions for privacy with no restrictions on the user's prior knowledge.

Section 4 delves deeper into the possibilistic model. For certain important cases, notably when the constraints on a user's prior knowledge are intersection-closed (i. e. not violated by the collusion of users), we give necessary and sufficient criteria for testing possibilistic privacy, which reduce the complexity of this problem.

Sections 5 and 6 focus on the more complex probabilistic model, over the set $\{0, 1\}^n$ of Boolean vectors that represent subsets of database records. Section 5 studies two probabilistic prior knowledge constraints: bit-wise independence (product distributions) and log-supermodularity. The bit-wise independence constraint was used also in [21] by Miklau and Suciu, so our work can be viewed as an extension of theirs. Log-supermodularity is chosen to provide a "middle ground" between bit-wise independence and the unconstrained prior knowledge. We give simple combinatorial necessary criteria and sufficient criteria for privacy under the log-supermodular and the product distribution constraints.

In Section 6, we study more general families Π of distributions over $\{0, 1\}^n$ that can be described by the intersection of a finite number of polynomial inequalities in a finite number of real-valued variables. We prove that even for certain very restricted Π , deciding whether a set $B \subseteq \{0, 1\}^n$ violates the privacy of a set $A \subseteq \{0, 1\}^n$ with respect to distributions in Π cannot be done in polynomial time, unless $P = NP$.

We overcome this negative result in two ways. First, using some deep results from algebraic geometry, we show that in certain interesting cases, such as when Π is the family of product distributions, there are provably efficient algorithms for deciding privacy. Second, we describe the sum-of-squares heuristic for deciding privacy for any Π , which has been implemented and works remarkably well in practice.

2. WORLDS AND AGENTS

Epistemology, the study of knowledge, has a long and honorable tradition in philosophy, starting with the early Greek philosophers. Philosophers were concerned with questions such as "What does it mean to say that someone knows something?" In the 1950's and 1960's [17, 19, 33] the focus shifted more to developing an *epistemic logic*, a logic of knowledge, and trying to capture the inherent properties of knowledge. Here there is a set Ω of *possible worlds*, one of which is the "real world" ω^* . An agent's *knowledge* is a set $S \subseteq \Omega$ of worlds that the agent considers possible. Since we are modeling *knowledge* rather than *belief*, we require that $\omega^* \in S$. If F is a (possible) fact, and $A \subseteq \Omega$ is the set of possible worlds where F is true, then we say that the agent *knows* F iff $S \subseteq A$.

More recently, researchers in such diverse fields as economics, linguistics, artificial intelligence, and theoretical computer science

have become interested in reasoning about knowledge [13]. The focus of attention has shifted to pragmatic concerns about the relationship between knowledge and action. That is our focus: the effect of an action, such as the disclosure of certain information, on the knowledge of an agent.

Worlds Let Ω be a finite set of all possible databases. We shall call a database $\omega \in \Omega$ a world, and the entire Ω the set of all possible worlds. The actual world, denoted by ω^* , represents the real database. Every property of the database, or assertion about its contents, can be formulated as " $\omega^* \in A$ " where $A \subseteq \Omega$ is the set of all databases that satisfy the property. A subset $A \subseteq \Omega$ that contains ω^* shall be called a *knowledge set*.

Agents We shall think of database users as *agents* who know something about the worlds in Ω and who try to figure out which $\omega \in \Omega$ is the actual world ω^* . An agent's knowledge can be modelled in different ways; we shall consider two approaches. In a *possibilistic* agent, knowledge is represented by a set $S \subseteq \Omega$ that contains exactly all the worlds this agent considers possible. In particular, $\omega^* \in S$. Here every world is either possible or not, with no ranking or score assigned. In a *probabilistic* agent, knowledge is represented by a probability distribution $P : \Omega \rightarrow \mathbb{R}_+$ that assigns a nonnegative weight $P(\omega)$ to every world. We denote the sum $\sum_{\omega \in A} P(\omega)$ by $P[A]$, requiring that $P[\Omega] = 1$ and $P(\omega^*) > 0$.

We say that a possibilistic agent with knowledge S *knows* a property $A \subseteq \Omega$ when $S \subseteq A$. We say that A is *possible* for this agent when $S \cap A \neq \emptyset$, i. e. when the agent does not know $\Omega - A$. For a probabilistic agent with distribution P , to *know* A means to have $P[A] = 1$, and to consider A possible means to have $P[A] > 0$.

A function Q that maps Ω to another set shall be called a *query*; if its range is $\{0, 1\}$ then Q is a *Boolean* query. For a given actual world ω^* , each query Q corresponds to the knowledge set associated with the query's "actual" output: $\{\omega \in \Omega \mid Q(\omega) = Q(\omega^*)\}$.

The Auditor There is a special "meta-agent" called the *auditor* whose task is to analyse the queries disclosed to the users and determine which of these disclosures could breach privacy. The auditor may or may not have complete information about the actual world ω^* . For example, if the query disclosure occurred several years ago, the record update logs may provide only a partial description of the database state at that moment. Even more importantly, the auditor does not know what the user's knowledge of the database was at the disclosure time. We characterize the auditor's knowledge by specifying which pairs of a database ω and the user's knowledge S (or P) the auditor considers possible. Let us formally define the auditor's knowledge about a user:

Definition 2.1. (Possibilistic case) A *possibilistic knowledge world* is a pair (ω, S) , where ω is a world and S is a knowledge set, which satisfies $\omega \in S \subseteq \Omega$. The set of all possibilistic knowledge worlds shall be denoted as

$$\Omega_{\text{poss}} := \{(\omega, S) \mid \omega \in S \subseteq \Omega\}$$

Ω_{poss} can be viewed as an extension of Ω . For a given user whose knowledge is $S^* \subseteq \Omega$, the pair $(\omega^*, S^*) \in \Omega_{\text{poss}}$ is called the *actual knowledge world*. The auditor's knowledge about the user is defined as a non-empty set $K \subseteq \Omega_{\text{poss}}$ of knowledge worlds, which must include the actual knowledge world. We refer to K as a *second-level knowledge set*.

Our knowledge worlds are similar to the 2-worlds of [14], except that the 2-worlds of [14] would deal not only with the knowledge that the user has of the world, but also with the knowledge that the auditor has of the world. Also, our second-level knowledge sets are similar to the 3-worlds of [14], except that the 3-worlds of [14] would deal not only with the knowledge that the auditor has about the user's knowledge of the world, but also with the knowledge that the user has about the auditor's knowledge of the world.

²Note that if Bob proactively tells Alice "If I am HIV-positive, then I had blood transfusions," a privacy breach of A may occur, because Alice may learn more than just B .

Definition 2.2. (Probabilistic case) A *probabilistic knowledge world* is a pair (ω, P) where P is a probability distribution over Ω such that $P(\omega) > 0$. The set of all probabilistic knowledge worlds shall be denoted as

$$\Omega_{\text{prob}} := \{(\omega, P) \mid P \text{ is a distribution, } P(\omega) > 0\}.$$

The actual knowledge world $(\omega^*, P^*) \in \Omega_{\text{prob}}$ and the auditor's second-level knowledge set $K \subseteq \Omega_{\text{prob}}$ are defined analogously to the possibilistic case.

Remark 2.3. The requirement of $\omega \in S$ for every pair $(\omega, S) \in \Omega_{\text{poss}}$ and of $P(\omega) > 0$ for every pair $(\omega, P) \in \Omega_{\text{prob}}$ represent our assumption that every agent considers the actual world possible. All pairs that violate this assumption are excluded as inconsistent. Note that a probabilistic pair (ω, P) is consistent iff the possibilistic pair $(\omega, \text{supp}(P))$ is consistent, where $\text{supp}(P) := \{\omega \mid P(\omega) > 0\}$.

Remark 2.4. In practice, it may be computationally infeasible to precisely characterize the auditor's second-level knowledge and to use this precisely characterized knowledge in the privacy definitions. Instead, the auditor makes assumptions about the database and the user's knowledge by placing constraints on the possible pairs (ω, S) or (ω, P) . These assumptions and constraints are also represented by a second-level knowledge set, which must contain the auditor's precise knowledge set as a subset. From now on, when we talk about the auditor's knowledge set, we mean a superset of the actual knowledge set, unless stated otherwise.

Definitions 2.1 and 2.2 allow us to consider an auditor whose assumptions about the user's knowledge depend on the contents of the database. For example, the auditor may assume that, if the hospital database contains record "Bob's doctor is Alice," then Alice knows Bob's HIV status, but if there is no such record, then Alice may or may not know it. However, in many situations we can separate the auditor's knowledge about the database from the auditor's assumptions about the user. We do so by specifying two sets:

1. A non-empty set $C \subseteq \Omega$ that consists of all databases the auditor considers possible, with $\omega^* \in C$;
2. A family Σ of subsets of Ω and/or a family Π of probability distributions over Ω . The possibilistic agent's knowledge has to belong to Σ , the probabilistic agent's knowledge has to belong to Π .

If the auditor knows the actual database exactly, e.g. by reconstructing its state from the update logs, then $C = \{\omega^*\}$; if the auditor has no information about the database or is unwilling to take advantage of it, then $C = \Omega$. Some choices for Σ and Π will be discussed in the subsequent sections.

When we say that the auditor's knowledge is represented by C and Σ described above, we mean that all knowledge worlds (ω, S) with $\omega \in C$ and $S \in \Sigma$, and none other, are considered possible by the auditor. However, in most cases the auditor's second-level knowledge set cannot be the Cartesian product $C \times \Sigma$, because it contains inconsistent (ω, S) pairs (see Remark 2.3). The same is true in the probabilistic case, for C and Π . Let us then define a product operation that excludes all inconsistent pairs:

Definition 2.5. The *product* of a set $C \subseteq \Omega$ and a family Σ of subsets of Ω (a family Π of probability distributions over Ω) is a second-level knowledge set $C \otimes \Sigma$ ($C \otimes \Pi$) defined by

$$\begin{aligned} C \otimes \Sigma &:= \{(\omega, S) \in C \times \Sigma \mid \omega \in S\} = (C \times \Sigma) \cap \Omega_{\text{poss}} \\ C \otimes \Pi &:= \{(\omega, P) \in C \times \Pi \mid P(\omega) > 0\} = (C \times \Pi) \cap \Omega_{\text{prob}} \end{aligned}$$

We call the pair (C, Σ) or (C, Π) *consistent* if their product $C \otimes \Sigma$ or $C \otimes \Pi$ is non-empty, because \emptyset is not a valid second-level knowledge set.

3. PRIVACY OF KNOWLEDGE

This section introduces the definition of privacy for the possibilistic and the probabilistic knowledge models. Let $A, B \subseteq \Omega$ be two arbitrary non-empty subsets of Ω ; as a shorthand, write $\bar{A} = \Omega - A$ and $AB = A \cap B$. Sets A and B correspond to two Boolean queries on the database ω^* ; e.g. query A returns "true" iff $\omega^* \in A$ and "false" otherwise.

We shall study the following question: When could the disclosure of B violate the privacy of A ? In our model, a positive result of query A is considered private and needs protection, whereas a negative result (that asserts \bar{A}) is not protected. Neither the user nor the auditor are assumed to know if A is true, and A may actually be false. On the other hand, B represents the disclosed fact, and therefore B has to be true. The auditor knows that B is true; the user transitions from not knowing B to knowing B .

Conceptually, we say that property A is private, given the disclosure of property B , if the user could not gain confidence in A by learning B . Below we shall make this notion precise for the two knowledge models, possibilistic and probabilistic. From now on, we shall use pronoun "he" for the user and "she" for the auditor.

3.1 Possibilistic Privacy

Let us suppose first that the auditor knows everything: the actual database ω^* such that $\omega^* \in B$, and the actual knowledge set S^* of the user at the time of the disclosure. In the possibilistic model, the user may have only two "grades of confidence" in property A : he either knows A ($S^* \subseteq A$), or he does not ($S^* \not\subseteq A$). The user gains confidence iff he does not know A before learning B (i.e. $S^* \not\subseteq A$) and knows A after learning B (i.e. $S^* \cap B \subseteq A$). Therefore, the privacy of A is preserved iff $\neg(S^* \not\subseteq A \ \& \ S^* \cap B \subseteq A)$, or equivalently, iff

$$S^* \cap B \subseteq A \Rightarrow S^* \subseteq A. \quad (4)$$

Now, suppose that the auditor does not know ω^* and S^* precisely, but has a second-level knowledge set $K \subseteq \Omega_{\text{poss}}$ such that $(\omega^*, S^*) \in K$. Then the auditor makes sure that A is private given B by checking condition (4) for all pairs in K . Before doing so, the auditor must discard from K all pairs (ω, S) such that $\omega \notin B$, because they are inconsistent with the disclosure of B . We arrive at the following possibilistic privacy definition:

Definition 3.1. Set $A \subseteq \Omega$ is called *K-private* given the disclosure of set $B \subseteq \Omega$, for $K \subseteq \Omega_{\text{poss}}$, iff

$$\forall (\omega, S) \in K : (\omega \in B \ \& \ S \cap B \subseteq A) \Rightarrow S \subseteq A. \quad (5)$$

We denote this predicate by $\text{Safe}_K(A, B)$.

Remark 3.2. It is easy to see from (5) that $\text{Safe}_K(A, B)$ and $K' \subseteq K$ imply $\text{Safe}_{K'}(A, B)$. Therefore, the auditor may assume less than she actually knows, i.e. consider more knowledge worlds possible, and still catch all privacy violations, at the expense of restricting more queries.

When the auditor wants to separate her knowledge about the database from her assumptions about the user's knowledge, she represents her second-level knowledge set K as a product $C \otimes \Sigma$, where $C \subseteq \Omega$ and Σ is a family of subsets of Ω . In this case we shall use the term " (C, Σ) -private" and the notation $\text{Safe}_{C, \Sigma}(A, B)$, which is defined as $\text{Safe}_{C \otimes \Sigma}(A, B)$; let us also use $\mathcal{P}(\Omega)$ to denote the power set of Ω .

PROPOSITION 3.3. For a consistent pair (C, Σ) such that $C \subseteq \Omega$ and $\Sigma \subseteq \mathcal{P}(\Omega)$, the privacy predicate $\text{Safe}_{C, \Sigma}(A, B)$ can be equivalently defined as follows: (denoting $S \cap B \cap C$ as SBC)

$$\forall S \in \Sigma : (SBC \neq \emptyset \ \& \ SB \subseteq A) \Rightarrow S \subseteq A. \quad (6)$$

3.2 Probabilistic Privacy

Once again, suppose first that the auditor knows the actual database $\omega^* \in B$ and the actual probability distribution P^* that represents the user's knowledge prior to the disclosure. As opposed to Section 3.1, in the probabilistic model the user has a continuum of "grades of confidence" in A , measured by $P^*[A]$. The user gains confidence iff his *prior* probability of A before learning B , which is $P^*[A]$, is strictly smaller than his *posterior* probability of A after B is disclosed, which is $P^*[A|B]$. Therefore, the privacy of A is preserved iff

$$P^*[A|B] \leq P^*[A]. \quad (7)$$

The conditional probability $P^*[A|B]$ is well-defined since $P^*[B] \geq P^*(\omega^*) > 0$.

When the auditor does not know ω^* and P^* , but has a second-level knowledge set $K \subseteq \Omega_{\text{prob}}$ such that $(\omega^*, P^*) \in K$, she has to check inequality (7) for all possible pairs (ω, P) in K . Before doing so, she must discard all pairs (ω, P) such that $\omega \notin B$. We obtain the following probabilistic privacy definition:

Definition 3.4. Set $A \subseteq \Omega$ is called K -private given the disclosure of set $B \subseteq \Omega$, for $K \subseteq \Omega_{\text{prob}}$, iff

$$\forall (\omega, P) \in K : \omega \in B \Rightarrow P[A|B] \leq P[A]. \quad (8)$$

As before, we denote this predicate by $\text{Safe}_K(A, B)$.

Remark 3.5. In the probabilistic case, too, $\text{Safe}_K(A, B)$ and $K' \subseteq K$ imply $\text{Safe}_{K'}(A, B)$. Thus, Remark 3.2 applies here.

When the auditor's knowledge can be represented as a product $C \otimes \Pi$ for some $C \subseteq \Omega$ and some family Π of probability distributions over Ω , we shall use the term " (C, Π) -private" and the notation $\text{Safe}_{C, \Pi}(A, B)$, which is defined as $\text{Safe}_{C \otimes \Pi}(A, B)$. In this case the following proposition can be used:

PROPOSITION 3.6. For a consistent pair (C, Π) where $C \subseteq \Omega$ and Π is a family of distributions over Ω , the privacy predicate $\text{Safe}_{C, \Pi}(A, B)$ can be equivalently defined as follows:

$$\forall P \in \Pi : P[BC] > 0 \Rightarrow P[AB] \leq P[A]P[B]. \quad (9)$$

In fact, the definition of privacy given by (9) can be further simplified, for many families Π that occur in practice:

Definition 3.7. We shall call a family Π ω -liftable for $\omega \in \Omega$ when $\forall P \in \Pi$ such that $P(\omega) = 0$ it satisfies the condition

$$\forall \varepsilon > 0 \exists P' \in \Pi : P'(\omega) > 0 \ \& \ \|P - P'\|_\infty < \varepsilon. \quad (10)$$

Family Π is called S -liftable for a set $S \subseteq \Omega$ iff Π is ω -liftable for $\forall \omega \in S$. The norm $\|P - P'\|_\infty := \max_{\omega \in \Omega} |P(\omega) - P'(\omega)|$.

PROPOSITION 3.8. For a consistent pair (C, Π) such that family Π is C -liftable, and given $BC \neq \emptyset$ (since $\omega^* \in BC$), predicate $\text{Safe}_{C, \Pi}(A, B)$ is equivalent to $\text{Safe}_\Pi(A, B)$ defined as

$$\text{Safe}_\Pi(A, B) \stackrel{\text{def}}{\iff} \forall P \in \Pi : P[AB] \leq P[A]P[B]. \quad (11)$$

3.3 Knowledge Acquisition

An agent (a database user) modifies his knowledge when he receives a disclosed query result. The disclosed property is a knowledge set $B \subseteq \Omega$, telling the agent that every world in $\Omega - B$ is impossible. We model the agent's acquisition of B as follows. A possibilistic agent with prior knowledge $S \subseteq \Omega$, upon receiving B such that $S \cap B \neq \emptyset$ (because $\omega^* \in S \cap B$), ends up with posterior knowledge $S \cap B$. A probabilistic agent with prior distribution $P : \Omega \rightarrow \mathbb{R}_+$, upon receiving B such that $P[B] \geq P(\omega^*) > 0$, ends up with posterior distribution $P(\cdot | B)$ defined by

$$P(\omega | B) = \begin{cases} P(\omega)/P[B], & \omega \in B \\ 0, & \omega \in \Omega - B \end{cases}$$

The acquisition of B_1 followed by B_2 is equivalent to the acquisition of $B_1 B_2 = B_1 \cap B_2$. When the auditor's second-level knowledge set K represents her assumption about the user's knowledge, rather than her knowledge of the user's knowledge (see Remark 2.4), she may want to require that K remains a valid assumption after each disclosure. This property is formalized below:

Definition 3.9. Let K be a second-level knowledge set, which may be possibilistic ($K \subseteq \Omega_{\text{poss}}$) or probabilistic ($K \subseteq \Omega_{\text{prob}}$). A set $B \subseteq \Omega$ is called K -preserving when for all $(\omega, S) \in K$ or $(\omega, P) \in K$ such that $\omega \in B$ we have $(\omega, S \cap B) \in K$ or $(\omega, P(\cdot | B)) \in K$.

Suppose that knowledge sets B_1 and B_2 are individually safe to disclose, while protecting the privacy of A , to an agent whose knowledge satisfies the constraints defined by K . If, after B_1 is disclosed, the updated agent's knowledge still satisfies the constraints, then it is safe to disclose B_2 too. Thus, it is safe to disclose both sets at once—as long as at least one of them preserves the constraints:

PROPOSITION 3.10. For every second-level knowledge set K , possibilistic or probabilistic, we have:

1. B_1 and B_2 are K -preserving $\Rightarrow B_1 \cap B_2$ is K -preserving;
2. If $\text{Safe}_K(A, B_1)$ and $\text{Safe}_K(A, B_2)$ and if at least one of B_1, B_2 is K -preserving, then $\text{Safe}_K(A, B_1 \cap B_2)$.

3.4 Unrestricted Prior Knowledge

What is the characterization of privacy when the auditor knows nothing? More formally, which knowledge sets A and B satisfy K -privacy for $K = \Omega_{\text{poss}} = \Omega \otimes \mathcal{P}(\Omega)$ and for $K = \Omega_{\text{prob}} = \Omega \otimes \mathcal{P}^{\text{prob}}(\Omega)$, where $\mathcal{P}^{\text{prob}}(\Omega)$ is the set of all probability distributions over Ω ? Also, what is the answer to this question if the auditor has complete information about the actual world ω^* , but knows nothing about the user's knowledge, i. e. for $K = \{\omega^*\} \otimes \mathcal{P}(\Omega)$ and for $K = \{\omega^*\} \otimes \mathcal{P}^{\text{prob}}(\Omega)$? Here is a theorem that answers these questions:

THEOREM 3.11. For all sets $A, B \subseteq \Omega$ and for all $\omega^* \in B$ the following four conditions are equivalent:

1. Either $A \cap B = \emptyset$, or $A \cup B = \Omega$;
2. $\text{Safe}_K(A, B)$ for $K = \Omega_{\text{poss}}$;
3. $\text{Safe}_K(A, B)$ for $K = \Omega_{\text{prob}}$;
4. $\text{Safe}_K(A, B)$ for $K = \{\omega^*\} \otimes \mathcal{P}^{\text{prob}}(\Omega)$.

Also, the following two conditions are equivalent (again $\omega^* \in B$):

1. $A \cap B = \emptyset$, or $A \cup B = \Omega$, or $\omega^* \in B - A$;
2. $\text{Safe}_K(A, B)$ for $K = \{\omega^*\} \otimes \mathcal{P}(\Omega)$.

Remark 3.12. In the auditing practice, the interesting case is $\omega^* \in A \cap B$, i. e. when the protected and the disclosed properties are both true. In this case $A \cap B \neq \emptyset$ and $\omega^* \notin B - A$. Unconditional privacy can thus be tested by checking whether $A \cup B = \Omega$, i. e. whether "A or B" is always true.

4. POSSIBILISTIC CASE

In this section, we shall focus exclusively on the possibilistic case; thus $K \subseteq \Omega_{\text{poss}}$. Proposition 4.1 below gives a necessary and sufficient condition for K -preserving sets B to satisfy the privacy predicate $\text{Safe}_K(A, B)$, for a given and fixed set A . It associates every world $\omega \in A$ with a "safety margin" $\beta(\omega) \subseteq \Omega - A$ which depends only on ω, A and K . Given B , the condition verifies whether every $\omega \in A$ occurs in B together with its "safety margin," or does not occur in B at all. The "safety margin" ensures that this ω will not reveal A to the agent, no matter what prior knowledge $S \in \pi_2(K)$ the agent might have. (By π_i we denote the projection operation.)

PROPOSITION 4.1. Let $K \subseteq \Omega_{\text{poss}}$ be an arbitrary second-level knowledge set, and let $A \subseteq \Omega$. There exists a function $\beta: A \rightarrow \mathcal{P}(\Omega - A)$ such that $\forall B \subseteq \Omega$

$$(\forall \omega \in AB : \beta(\omega) \subseteq B) \Rightarrow \text{Safe}_K(A, B), \quad (12)$$

and if B is K -preserving, then the converse holds:

$$\text{Safe}_K(A, B) \Rightarrow (\forall \omega \in AB : \beta(\omega) \subseteq B). \quad (13)$$

Remark 4.2. In Proposition 4.1, the condition of B being K -preserving is essential for the converse implication (13). Indeed, let $\Omega = \{1, 2, 3\}$, $K = \Omega \otimes \{\Omega\}$, and $A = \{3\}$. Then both $B_1 = \{1, 3\}$ and $B_2 = \{2, 3\}$ protect the K -privacy of A , yet $B_1 \cap B_2 = \{3\}$ does not. So, there is no suitable value for $\beta(3)$. However, see Corollary 4.14 for more on this subject.

The characterization in Proposition 4.1 could be quite useful for auditing a lot of properties B_1, B_2, \dots, B_N disclosed over a period of time, using the same audit query A . Given A , the auditor would compute the mapping β once, and use it to test every B_i . This comment applies to Section 4.1 as well.

4.1 Intersection-Closed Knowledge

Motivation. When two or more possibilistic agents collude, i. e. join forces in attacking protected information, their knowledge sets intersect: they jointly consider a world possible if and only if none of them has ruled it out. Therefore, if the auditor wants to account for potential collusions, she must consider knowledge world $(\omega, S_1 \cap S_2)$ possible whenever she considers both (ω, S_1) and (ω, S_2) possible. This motivates the following definition:

Definition 4.3. A second-level knowledge set $K \subseteq \Omega_{\text{poss}}$ is *intersection-closed*, or \cap -closed for short, iff $\forall (\omega, S_1) \in K$ and $\forall (\omega, S_2) \in K$ we have $(\omega, S_1 \cap S_2) \in K$. Note that we intersect the user's knowledge sets (ω, S_1) and (ω, S_2) only when they are paired with the same world ω .

One way to obtain a second-level knowledge set $K \subseteq \Omega_{\text{poss}}$ that is \cap -closed is by taking an \cap -closed family Σ of subsets of Ω (such that $\forall S_1, S_2 \in \Sigma : S_1 \cap S_2 \in \Sigma$) and computing the product $K = C \otimes \Sigma$ with some knowledge set C .

Intervals. When the auditor's knowledge is \cap -closed, the notion of an "interval" between two worlds becomes central in characterizing the privacy relation:

Definition 4.4. Let $K \subseteq \Omega_{\text{poss}}$ be \cap -closed, and let $\omega_1, \omega_2 \in \Omega$ be two worlds such that

$$\omega_1 \in \pi_1(K), \quad \omega_2 \in \bigcup \{S \mid (\omega_1, S) \in K\}. \quad (14)$$

The K -interval from ω_1 to ω_2 , denoted by $I_K(\omega_1, \omega_2)$, is the smallest set S such that $(\omega_1, S) \in K$ and $\omega_2 \in S$, or equivalently:

$$I_K(\omega_1, \omega_2) := \bigcap \{S \mid (\omega_1, S) \in K, \omega_2 \in S\}.$$

If the worlds ω_1, ω_2 do not satisfy conditions (14), we shall say that interval $I_K(\omega_1, \omega_2)$ does not exist.

The following proposition shows that we need to know only the intervals in order to check whether or not $\text{Safe}_K(A, B)$ holds:

PROPOSITION 4.5. For an \cap -closed set $K \subseteq \Omega_{\text{poss}}$ and for all $A, B \subseteq \Omega$, we have $\text{Safe}_K(A, B)$ if and only if

$$\begin{aligned} \forall I_K(\omega_1, \omega_2) : \quad & \omega_1 \in AB \ \& \ \omega_2 \notin A \\ & \Rightarrow I_K(\omega_1, \omega_2) \cap (B - A) \neq \emptyset. \end{aligned} \quad (15)$$

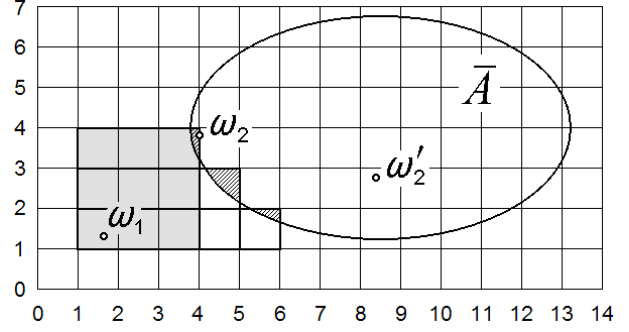


Figure 1: An example of an \cap -closed $K \subseteq \Omega_{\text{poss}}$ where the worlds are the pixels inside the 14×7 rectangle (such as ω_1 , ω_2 and ω'_2), and the permitted user's knowledge sets are the integer sub-rectangles (i. e. composed of whole squares). Set \bar{A} is the complement of the privacy-sensitive knowledge set. See Example 4.9 for details.

Remark 4.6. As implied by Proposition 4.5, there is no need to store the entire \cap -closed second-level knowledge set K (which could require $|\Omega| \cdot 2^{|\Omega|}$ bits of data) in order to test the possibilistic privacy. It is sufficient to store one set $I_K(\omega_1, \omega_2) \subseteq \Omega$, or the fact of its non-existence, for each pair $(\omega_1, \omega_2) \in \Omega \times \Omega$, i. e. at most $|\Omega|^3$ bits of data.

Minimal intervals. In fact, in Proposition 4.5 we do not even have to check all intervals; it is enough to consider just "minimal" intervals defined as follows:

Definition 4.7. For an \cap -closed second-level knowledge set $K \subseteq \Omega_{\text{poss}}$, for a world $\omega_1 \in \Omega$ and for a set $X \subseteq \Omega$ not containing ω_1 , an interval $I_K(\omega_1, \omega_2)$ is called a *minimal K -interval from ω_1 to X* iff $\omega_2 \in X$ and

$$\forall \omega'_2 \in X \cap I_K(\omega_1, \omega_2) : I_K(\omega_1, \omega'_2) = I_K(\omega_1, \omega_2).$$

PROPOSITION 4.8. For an \cap -closed set $K \subseteq \Omega_{\text{poss}}$ and for $\forall A, B \subseteq \Omega$, we have $\text{Safe}_K(A, B)$ if and only if the formula (15) holds over all intervals $I_K(\omega_1, \omega_2)$ that are minimal from a world $\omega_1 \in AB$ to the set $\Omega - A$.

Example 4.9. Let Ω be an area of the plane that is bounded by a rectangle and discretized into pixels to ensure finiteness (the area within the 14×7 rectangle on Figure 1). Let the worlds be the pixels. Consider an auditor who does not know the actual database ω^* and who assumes that the user's prior knowledge set $S \in \Sigma$ is an integer rectangle, i. e. a rectangle whose four corners have integer coordinates (corresponding to the vertical and horizontal lines in the picture). The family Σ of integer rectangles is \cap -closed, and so is the auditor's second-level knowledge set $K = \Omega \otimes \Sigma$.

Given $\omega_1, \omega_2 \in \Omega$, the interval $I_K(\omega_1, \omega_2)$ is the smallest integer rectangle that contains both ω_1 and ω_2 . For ω_1 and ω_2 in Figure 1, the interval $I_K(\omega_1, \omega_2)$ is the light-grey rectangle from point $(1, 1)$ to point $(4, 4)$; for ω_1 and ω'_2 , the interval $I_K(\omega_1, \omega'_2)$ is the rectangle from point $(1, 1)$ to point $(9, 3)$.

The interval $I_K(\omega_1, \omega_2)$ shown on the picture is one of the three minimal intervals from ω_1 to set \bar{A} (the area bounded by the ellipse). The other two minimal intervals are the rectangles $(1, 1) - (5, 3)$ and $(1, 1) - (6, 2)$. Every set S such that $(\omega_1, S) \in K$ and $S \not\subseteq \bar{A}$, e. g. the interval $I_K(\omega_1, \omega'_2)$, must contain at least one of the three minimal intervals, implying Proposition 4.8 for the case of the actual world $\omega^* = \omega_1$.

Interval-induced partitions of \bar{A} . Let us have a closer look at the minimal K -intervals from a given world $\omega_1 \in A$ to the set

$\bar{A} = \Omega - A$. For every $\omega_2 \in \bar{A}$, the interval $I_K(\omega_1, \omega_2)$, if it exists, is either minimal or not; if it is not minimal, then ω_2 cannot belong to any minimal interval from ω_1 to \bar{A} . Now, take some pair $\omega_2, \omega'_2 \in \bar{A}$ such that both $I_K(\omega_1, \omega_2)$ and $I_K(\omega_1, \omega'_2)$ are minimal. There are two possible situations:

1. $I_K(\omega_1, \omega_2) = I_K(\omega_1, \omega'_2)$, or
2. $I_K(\omega_1, \omega_2) \cap I_K(\omega_1, \omega'_2) \cap \bar{A} = \emptyset$.

Indeed, if $\exists \omega'_2 \in I_K(\omega_1, \omega_2) \cap I_K(\omega_1, \omega'_2) \cap \bar{A}$, then by Definition 4.7 the interval $I_K(\omega_1, \omega'_2)$ equals both of the minimal intervals, making them equal. We have thus shown the following

PROPOSITION 4.10. *Given an \cap -closed set $K \subseteq \Omega_{\text{poss}}$, a set $A \subseteq \Omega$, and a world $\omega_1 \in A$, the minimal K -intervals from ω_1 to \bar{A} partition set \bar{A} into disjoint equivalence classes*

$$\bar{A} = D_1 \cup D_2 \cup \dots \cup D_m \cup D'$$

where two worlds $\omega_2, \omega'_2 \in \bar{A}$ belong to the same class D_i iff they both belong to the same minimal interval, or (class D') when they both do not belong to any minimal interval.

Definition 4.11. In the assumptions and in the notation of Proposition 4.10, denote

$$\Delta_K(\bar{A}, \omega_1) := \{D_1, D_2, \dots, D_m\}.$$

In other words, $\Delta_K(\bar{A}, \omega_1)$ is the disjoint collection of all sets formed by intersecting \bar{A} with the minimal intervals from ω_1 to \bar{A} .

COROLLARY 4.12. *Given an \cap -closed set $K \subseteq \Omega_{\text{poss}}$, for all $A, B \subseteq \Omega$ we have $\text{Safe}_K(A, B)$ if and only if*

$$\forall \omega_1 \in AB, \forall D_i \in \Delta_K(\bar{A}, \omega_1) : B \cap D_i \neq \emptyset. \quad (16)$$

As Figure 1 illustrates for Example 4.9, the three minimal intervals from ω_1 to \bar{A} formed by integer rectangles $(1, 1) - (4, 4)$, $(1, 1) - (5, 3)$ and $(1, 1) - (6, 2)$ are disjoint inside \bar{A} . Their intersections with \bar{A} , shown hatched in Figure 1, constitute the collection $\Delta_K(\bar{A}, \omega_1)$. A disclosed set B is private, assuming $\omega^* = \omega_1$, iff B intersects each of these three intervals inside \bar{A} .

The case of all-singleton Δ_K 's. If set K satisfies the property defined next, privacy testing is simplified still further:

Definition 4.13. An \cap -closed set $K \subseteq \Omega_{\text{poss}}$ has *tight intervals* iff for every K -interval $I_K(\omega_1, \omega_2)$ we have

$$\forall \omega'_2 \in I_K(\omega_1, \omega_2) : \omega'_2 \neq \omega_2 \Rightarrow I_K(\omega_1, \omega'_2) \subsetneq I_K(\omega_1, \omega_2).$$

When K has tight intervals, every minimal interval $I_K(\omega_1, \omega_2)$ from $\omega_1 \in A$ to \bar{A} has *exactly one* of its elements in \bar{A} , namely ω_2 : $\bar{A} \cap I_K(\omega_1, \omega_2) = \{\omega_2\}$. Then all equivalence classes D_i in $\Delta_K(\bar{A}, \omega_1)$ are singletons, and Corollary 4.12 takes the form of Proposition 4.1:

COROLLARY 4.14. *Let $K \subseteq \Omega_{\text{poss}}$ be an \cap -closed set that has tight intervals, let $A \subseteq \Omega$. Then $\exists \beta : A \rightarrow \mathcal{P}(\Omega - A)$ given by*

$$\forall \omega_1 \in A : \beta(\omega_1) := \bigcup \Delta_K(\bar{A}, \omega_1)$$

such that $\forall B \subseteq \Omega$

$$\text{Safe}_K(A, B) \Leftrightarrow (\forall \omega_1 \in AB : \beta(\omega_1) \subseteq B).$$

Having tight intervals is essential for Corollary 4.14 to hold; see Remark 4.2 for a counterexample where an \cap -closed K does not have tight intervals.

5. MODULARITY ASSUMPTIONS FOR PROBABILISTIC KNOWLEDGE

In the previous section we clarified some general properties of possibilistic knowledge; now we turn to the more complex probabilistic case. Rather than studying arbitrary probabilistic knowledge families, here we shall focus on a few specific, yet important, families of distributions. Later, in Section 6, we present more sophisticated approaches that extend beyond these families.

From now on, we assume that $\Omega = \{0, 1\}^n$ for some fixed n . Let $\omega_1 \wedge \omega_2$ ($\omega_1 \vee \omega_2$, $\omega_1 \oplus \omega_2$) be the bit-wise ‘‘AND’’ (‘‘OR’’, ‘‘XOR’’), and define the partial order $\omega_1 \leq \omega_2$ to mean ‘‘ $\forall i = 1 \dots n : \omega_1[i] = 1 \Rightarrow \omega_2[i] = 1$.’’ A set $S \subseteq \Omega$ shall be called an *up-set* (a *down-set*) when $\forall \omega_1 \in S, \forall \omega_2 \geq \omega_1$ ($\forall \omega_2 \leq \omega_1$) we have $\omega_2 \in S$.

Definition 5.1. A probability distribution P over Ω is called *log-supermodular* (*log-submodular*)³ when the following holds:

$$\forall \omega_1, \omega_2 \in \Omega : P(\omega_1)P(\omega_2) \leq (\geq) P(\omega_1 \wedge \omega_2)P(\omega_1 \vee \omega_2)$$

The family of all log-supermodular distributions shall be denoted by Π_m^+ , the family of all log-submodular distributions by Π_m^- .

A distribution P is called a *product distribution* if it makes every coordinate independent. Every product distribution corresponds to a vector (p_1, \dots, p_n) of Bernoulli probabilities, each $p_i \in [0, 1]$, such that

$$\forall \omega \in \{0, 1\}^n : P(\omega) = \prod_{i=1}^n p_i^{\omega[i]} \cdot (1 - p_i)^{1 - \omega[i]} \quad (17)$$

The family of all product distributions shall be denoted by Π_m^0 . It is easy to show that $\Pi_m^0 = \Pi_m^- \cap \Pi_m^+$ [20]. In fact, P is a product distribution if and only if

$$\forall \omega_1, \omega_2 \in \Omega : P(\omega_1)P(\omega_2) = P(\omega_1 \wedge \omega_2)P(\omega_1 \vee \omega_2) \quad (18)$$

Supermodular and submodular functions occur often in mathematics and have been extensively studied [15, 20]. Our goal in considering these assumptions was to substantially relax bit-wise independence while staying away from the unconstrained case. Besides that, the log-supermodular assumption (as implied by Theorem 5.3 below) describes situations where no negative correlations are permitted between positive events—something we might expect from knowledge about, say, HIV incidence among humans.

PROPOSITION 5.2 (Π_m^+ SAFETY: NECESSARY CRITERION).

For all $A, B \subseteq \Omega = \{0, 1\}^n$, we have:

$$\text{Safe}_{\Pi_m^+}(A, B) \Rightarrow \forall \omega_1 \in AB, \forall \omega_2 \in \bar{A}\bar{B} : \quad (19)$$

$$\left(\begin{array}{l} \omega_1 \wedge \omega_2 \in A - B \\ \omega_1 \vee \omega_2 \in B - A \end{array} \right) \text{ or } \left(\begin{array}{l} \omega_1 \wedge \omega_2 \in B - A \\ \omega_1 \vee \omega_2 \in A - B \end{array} \right)$$

Our sufficient criterion for Π_m^+ -safety has a very similar form, and relies on the following well-known theorem [3] (see also [6], §19):

THEOREM 5.3 (FOUR FUNCTIONS THEOREM). *Let L be a distributive lattice, and let $\alpha, \beta, \gamma, \delta : L \rightarrow \mathbb{R}_+$. For all $A, B \subseteq L$ denote $f[A] = \sum_{a \in A} f(a)$, $A \vee B = \{a \vee b \mid a \in A, b \in B\}$, and $A \wedge B = \{a \wedge b \mid a \in A, b \in B\}$. Then the inequality*

$$\alpha[A] \cdot \beta[B] \leq \gamma[A \vee B] \cdot \delta[A \wedge B]$$

holds for all subsets $A, B \subseteq L$ if and only if it holds for one-element subsets, i. e. iff

$$\alpha(a) \cdot \beta(b) \leq \gamma(a \vee b) \cdot \delta(a \wedge b)$$

for all elements $a, b \in L$.

³The ‘‘log-’’ means that supermodularity is multiplicative, not additive. The subscript ‘‘m’’ in Π_m^- , Π_m^+ etc. means ‘‘modular.’’

PROPOSITION 5.4 (Π_m^+ SAFETY: SUFFICIENT CRITERION).

For all $A, B \subseteq \Omega = \{0, 1\}^n$, either one of the two conditions below is sufficient to establish $\text{Safe}_{\Pi_m^+}(A, B)$:

- $AB \wedge \bar{A}\bar{B} \subseteq A - B$ and $AB \vee \bar{A}\bar{B} \subseteq B - A$;
- $AB \vee \bar{A}\bar{B} \subseteq A - B$ and $AB \wedge \bar{A}\bar{B} \subseteq B - A$.

COROLLARY 5.5. If A is an up-set and B is a down-set (or vice versa), then $\text{Safe}_{\Pi_m^+}(A, B)$.

Remark 5.6. Thus, if the user's prior knowledge is assumed to be in Π_m^+ , a "no" answer to a monotone Boolean query always preserves the privacy of a "yes" answer to another monotone Boolean query. Roughly speaking, it is OK to disclose a negative fact while protecting a positive fact. This observation is especially helpful when A and B are given by query language expressions, whose monotonicity is often obvious.

5.1 Product Distributions

In this section we shall study the problem of checking the privacy relation $\text{Safe}_{\Pi_m^0}(A, B)$ for sets $A, B \subseteq \Omega = \{0, 1\}^n$ over the family Π_m^0 of product distributions. The *independence* relation that holds iff $P[A]P[B] = P[AB]$ for all $P \in \Pi_m^0$, and which we denote by $A \perp_{\Pi_m^0} B$, has been studied by Miklau and Suciu in [21] who proved the following necessary and sufficient criterion:

THEOREM 5.7 (MIKLAU & SUCIU). For all $A, B \subseteq \Omega$, $A \perp_{\Pi_m^0} B$ if and only if sets A and B "share no critical coordinates," i.e. when coordinates $1, 2, \dots, n$ can be rearranged so that only $\omega[1], \omega[2], \dots, \omega[k]$ determine if $\omega \in A$, and only $\omega[k+1], \omega[k+2], \dots, \omega[k']$, $k' \leq n$, determine if $\omega \in B$.

Since $A \perp_{\Pi_m^0} B$ implies $\text{Safe}_{\Pi_m^0}(A, B)$, the Miklau-Suciu criterion is a sufficient criterion for our notion of privacy. It is not a necessary one, even for $n = 2$: we have $\text{Safe}_{\Pi_m^0}(X_1, \bar{X}_1 \cup X_2)$ but not $X_1 \perp_{\Pi_m^0} (\bar{X}_1 \cup X_2)$, where $X_i = \{\omega \in \Omega \mid \omega[i] = 1\}$.

Another sufficient criterion is given by Proposition 5.4, if we note that $\Pi_m^0 \subset \Pi_m^+$; it implies $\text{Safe}_{\Pi_m^0}(A, B)$ whenever A is an up-set and B is a down-set, or vice versa (Corollary 5.5). A little more generally, $\text{Safe}_{\Pi_m^0}(A, B)$ holds if there exists a mask vector $z \in \Omega$ such that $z \oplus A$ is an up-set and $z \oplus B$ is a down-set. Let us call this criterion the *monotonicity criterion*.

It turns out that both the Miklau-Suciu and the monotonicity criteria are special cases of another simple yet surprisingly strong sufficient criterion for $\text{Safe}_{\Pi_m^0}(A, B)$. This sufficient criterion shall be called the *cancellation criterion*, because its verification is equivalent to cancelling identical monomial terms in the algebraic expansion for the difference

$$P[A\bar{B}]P[\bar{A}B] - P[AB]P[\bar{A}\bar{B}] = P[A]P[B] - P[AB],$$

where P is a product distribution written as in (17). In order to formulate the criterion in combinatorial (rather than algebraic) terms, we need the following definition:

Definition 5.8. The *pairwise matching function* $\text{Match}(u, v)$ maps a pair (u, v) of vectors from $\Omega = \{0, 1\}^n$ to a single *match-vector* $w = \text{Match}(u, v)$ in $\{0, 1, *\}^n$ as follows:

$$\forall i = 1 \dots n: \quad w[i] = \begin{cases} u[i] & \text{if } u[i] = v[i]; \\ * & \text{if } u[i] \neq v[i]. \end{cases}$$

For example, pair $(01011, 01101)$ gets mapped into $01**1$. We say that $v \in \Omega$ *refines* a match-vector w when v can be obtained from w by replacing its every star with a 0 or a 1. For every match-vector w , define the following two sets:

$$\begin{aligned} \text{Box}(w) &:= \{v \in \Omega \mid v \text{ refines } w\}; \\ \text{Circ}(w) &:= \{(u, v) \in \Omega \times \Omega \mid \text{Match}(u, v) = w\}. \end{aligned}$$

Now we are ready to state the cancellation criterion, which is a sufficient criterion for $\text{Safe}_{\Pi_m^0}(A, B)$, and also state a necessary criterion of a similar form, for comparison:

PROPOSITION 5.9 (CANCELLATION CRITERION). For all $A, B \subseteq \Omega$, in order to establish $\text{Safe}_{\Pi_m^0}(A, B)$ it is sufficient to verify the following $\forall w \in \{0, 1, *\}^n$:

$$|\bar{A}\bar{B} \times \bar{A}B \cap \text{Circ}(w)| \geq |AB \times \bar{A}\bar{B} \cap \text{Circ}(w)|.$$

PROPOSITION 5.10 (A NECESSARY CRITERION). For all $A, B \subseteq \Omega$, if $\text{Safe}_{\Pi_m^0}(A, B)$ holds, then $\forall w \in \{0, 1, *\}^n$:

$$\begin{aligned} |A\bar{B} \cap \text{Box}(w)| \cdot |\bar{A}B \cap \text{Box}(w)| &\geq \\ &\geq |AB \cap \text{Box}(w)| \cdot |\bar{A}\bar{B} \cap \text{Box}(w)|. \end{aligned}$$

We hope that the combinatorial simplicity of the criterion given by Proposition 5.9 will allow highly scalable implementations that apply in real-life database auditing scenarios, where sets A and B are given via expressions in a query language. The theorem below justifies our interest in the cancellation criterion:

THEOREM 5.11. If sets A, B satisfy the Miklau-Suciu criterion or the monotonicity criterion, they also satisfy the cancellation criterion.

Remark 5.12. The cancellation criterion is only sufficient, but not necessary. Here is a pair of sets that satisfies $\text{Safe}_{\Pi_m^0}(A, B)$ and does not satisfy the criterion: $A = \{011, 100, 110, 111\}$ and $B = \{010, 101, 110, 111\}$. Specifically, for these sets we have $|A\bar{B} \times \bar{A}B \cap \text{Circ}(***)| = 0$ and $|AB \times \bar{A}\bar{B} \cap \text{Circ}(***)| = 2$.

6. GENERAL ALGEBRAIC APPROACHES

We use techniques from multivariate polynomial optimization to test safety with respect to certain families Π of prior distributions on an agent's knowledge. Recall that a set $A \subseteq \Omega$ is Π -safe given $B \subseteq \Omega$ when for all distributions $P \in \Pi$, we have $P[AB] \leq P[A] \cdot P[B]$. As in some previous sections, we identify the set Ω of possible worlds with the hypercube $\{0, 1\}^n$.

For each $x \in \{0, 1\}^n$, we create variables $p_x \in [0, 1]$. We consider those families Π containing distributions $(p_x)_{x \in \{0, 1\}^n}$ which can be described by the intersection of a finite number r of polynomial inequalities:

$$\begin{aligned} \alpha_1((p_x)_{x \in \{0, 1\}^n}) \geq 0, \dots, \alpha_r((p_x)_{x \in \{0, 1\}^n}) \geq 0, \\ \sum_{x \in \{0, 1\}^n} p_x = 1, \quad \forall x \ p_x \geq 0. \end{aligned}$$

We call such a family Π *algebraic*. For example, if we had the family of log-submodular distributions, then for all $x, y \in \{0, 1\}^n$, we would have the constraint $\alpha_{x,y} = p_x p_y - p_{x \wedge y} p_{x \vee y} \geq 0$. For the family of log-supermodular distributions, we would instead have $\alpha_{x,y} = p_{x \wedge y} p_{x \vee y} - p_x p_y \geq 0$. Finally, for the family of product distributions, we would have both $p_x p_y - p_{x \wedge y} p_{x \vee y} \geq 0$ and $p_{x \wedge y} p_{x \vee y} - p_x p_y \geq 0$.

For sets A and B , and a family of distributions Π , we define the set $K(A, B, \Pi)$ of distributions $(p_x)_{x \in \{0, 1\}^n}$ to be:

$$\begin{aligned} \sum_{w \in AB} p_w > \sum_{x \in A} p_x \sum_{y \in B} p_y \\ \alpha_1((p_x)_{x \in \{0, 1\}^n}) \geq 0, \dots, \alpha_r((p_x)_{x \in \{0, 1\}^n}) \geq 0 \\ \sum_{x \in \{0, 1\}^n} p_x = 1, \quad \forall x \ p_x \geq 0. \end{aligned}$$

The following proposition is an equivalent algebraic formulation of the fact that in order for $\text{Safe}_{\Pi}(A, B)$ to hold, there cannot be a single distribution $P \in \Pi$ for which $P[AB] > P[A] \cdot P[B]$.

PROPOSITION 6.1. $\text{Safe}_\Pi(A, B)$ iff the set $K(A, B, \Pi)$ is empty.

We are interested in algorithms that decide emptiness of $K(A, B, \Pi)$ in time polynomial or nearly polynomial in $N \stackrel{\text{def}}{=} 2^n$. Note that N does not need to be the number of possible worlds, but rather only the potentially much smaller number of possible *relevant worlds* in the desired application. For example, if the agent executes a combination of PROJECT and SELECT queries in SQL, he may be left only with a subset S of possible records with a small number of attributes and values for those attributes. In this case, the number N of possible relevant worlds could be very small, and algorithms for testing safety of additional queries on S which run in time polynomial or quasi-polynomial in N would be efficient.

As the following theorem shows, even when the number N of possible relevant worlds is small, we may need to restrict the class of distributions Π that we consider in order to efficiently test safety.

THEOREM 6.2. *If $P \neq NP$, there is an algebraic Π for which $r = \text{poly}(N)$, each α_i has degree at most 2, and for which deciding $\text{Safe}_\Pi(A, B)$ cannot be done in $\text{poly}(N)$ time.*

PROOF. (sketch) The main idea is a reduction from a restricted version of the decision problem of MAX-CUT. We carefully choose constraints defining the family Π so that given a graph G on t vertices, we can encode G into sets $A, B \subseteq \{0, 1\}^n$ so that the constraints defining Π together with the constraint $P[AB] > P[A] \cdot P[B]$ define a non-empty set $K(A, B, \Pi)$ iff the maximum cut size in G is sufficiently large. We need to suitably restrict the decision version of MAX-CUT so that this is possible. Here we require $N = \text{poly}(t)$. We defer the details of the proof to the full paper. \square

Despite this negative result, for certain interesting families Π we obtain efficient algorithms, as we now discuss.

6.1 Specific Distributions

We first obtain a necessary and sufficient condition for $A, B \subseteq \{0, 1\}^n$ to be safe with respect to the family Π of product distributions by providing a deterministic algorithm. Its running time is $N^{O(\lg \lg N)}$, which is essentially polynomial for all practical purposes. The key observation is that while $K(A, B, \Pi)$ is $N = 2^n$ -dimensional for general families of distributions, for product distributions it can be embedded into \mathbb{R}^n .

Indeed, it is easy to see that $K(A, B, \Pi)$ can be defined in variables $p_1, \dots, p_n \in \mathbb{R}$ constrained by $p_i(1-p_i) \geq 0$, and for which $P[AB] > P[A] \cdot P[B]$, where $P(\omega) = \prod_{i=1}^n p_i^{\omega[i]} \cdot (1-p_i)^{1-\omega[i]}$ for all $\omega \in \{0, 1\}^n$. We can write this with n variables and $n+1$ inequalities. We apply the following simplified form of Theorem 3 of Basu, Pollack, and Roy [4]:

THEOREM 6.3. *Given a set $K = \{\beta_1, \dots, \beta_r\}$ of r polynomials each of degree at most d in s variables with coefficients in \mathbb{R} , the problem of deciding whether there exist $X_1, \dots, X_s \in \mathbb{R}$ for which $\beta_1(X_1, \dots, X_s) \geq 0, \dots, \beta_r(X_1, \dots, X_s) \geq 0$, can be solved deterministically with $\tau(rd)^{O(k)}$ bit operations, where τ is the number of bits needed to describe a coefficient in β_1, \dots, β_r .*

We apply this theorem to the set $K = K(A, B, \Pi)$. From the program above it is easy to see that τ, r, d , and s are all linear in n , and so emptiness (and hence safety) for product distributions can be decided in $n^{O(n)} = N^{O(\lg \lg N)}$ time.

The algorithm of Basu, Pollack, and Roy uses sophisticated ideas from algebraic geometry over \mathbb{R} , and we cannot do it justice here. The general approach taken by such algorithms is to reduce a system of polynomial inequalities into a system of polynomial equalities by introducing slack variables, and then combining the multivariate polynomial equalities $p_i(x) = 0$ into a single equality

$q(x) \stackrel{\text{def}}{=} \sum_x p_i^2(x) = 0$. One finds the critical points of $q(x)$, that is, the set V_C of common zeros of its partial derivatives over the complex field \mathbb{C} . By perturbing $q(x)$ and applying Bézout's Theorem, one can show that $|V_C|$ is finite. Various approaches are used to find the subset $V_{\mathbb{R}}$ of V_C of real-valued points. Since $V_{\mathbb{R}}$ is finite, once it is found q is evaluated on each of its elements and the minimum value is taken. The main step is finding $V_{\mathbb{R}}$, and approaches based on Gröbner bases, resultant theory, and homotopy theory exist (see [25]). The algorithm of [4] may be practical. Indeed, a similar algorithm of Canny was implemented [7].

This approach generalizes to other algebraic families Π described by $\text{poly}(n)$ constraints and $O(n)$ variables. For instance, a family of distributions for which $p_x = p_y$ whenever the Hamming weight of x and y are equal is described by $n+1$ variables.

Even when the family Π of distributions requires N variables to describe, in certain cases we can obtain a polynomial-time algorithm for testing safety with respect to Π . Indeed, if the constraints α_i defining Π have degree at most 2 and there are only a constant number r of them, an algorithm in [16] shows how to decide emptiness of $K(A, B, \Pi)$ in $N^{O(r)}$ time. This algorithm makes black-box use of the earlier algorithm of Basu, Pollack, and Roy [4]. As an optimization, we note that if there are multiple linear equality constraints $L_i(X_1, \dots, X_s) = 0$, it is helpful to combine them into a single quadratic constraint $\sum_i L_i^2 = 0$. This is because the running time is exponential in the number of constraints.

6.2 Heuristics

For most families of distributions we will have to settle for a heuristic or an approximation for testing safety. If the program describing $K(A, B, \Pi)$ is multilinear (e.g., one can show this is the case for log-submodular and log-supermodular distributions), there are heuristics such as branch-and-bound or cutting-plane techniques. See page 2 of [9].

Here we describe the arguably most practical heuristic, the *sum-of-squares* heuristic, introduced in [30, 31, 24], which works even for systems that are not multilinear. This heuristic was implemented with great success in [25].

The problem of minimizing a degree- d multivariate polynomial f over a set $K \subseteq \mathbb{R}^s$ is equivalent to finding the maximum $\gamma \in \mathbb{R}$ for which $f(x) - \gamma \geq 0$ for all $x \in K$. Let $\mathcal{P}_+^d(K)$ be the set of all polynomials in $\mathbb{R}[x_1, \dots, x_s]$ of degree at most d which are non-negative on every point in K . Thus, our problem is to find the maximum $\gamma \in \mathbb{R}$ for which $f(x) - \gamma \in \mathcal{P}_+^d(K)$.

It is unknown how to optimize over $\mathcal{P}_+^d(K)$ efficiently, and so the following indirect route is taken. Define the set Σ^2 :

$$\Sigma^2 = \{f(x) \in \mathbb{R}[x_1, \dots, x_s] \mid \exists t, g_1(x), \dots, g_t(x) \in \mathbb{R}[x_1, \dots, x_s] \text{ s.t. } f(x) = \sum_{i=1}^t g_i(x)^2\}.$$

Notice that Σ^2 is a subset of non-negative polynomials, as every sum of squares of polynomials is non-negative. It turns out that Σ^2 is in fact a strict subset of the non-negative polynomials, as shown non-constructively by Hilbert, and constructively by Motzkin who provided the polynomial

$$M(x, y, z) = x^4y^2 + x^2y^4 + z^6 - 3x^2y^2z^2.$$

Motzkin showed $M(x, y, z)$ is non-negative on \mathbb{R}^3 , yet inexpressible as a sum of squares of polynomials. It turns out that every non-negative polynomial can be written as a sum of squares of rational functions (functions of the form $g_i(x)/h_i(x)$ for polynomials g_i and h_i), which was Hilbert's 17th problem, solved by Artin in 1927. While Σ^2 fails to capture all non-negative polynomials, the following proposition is a compelling reason for studying it. The proposition is proven using semidefinite programming.

PROPOSITION 6.4. *For $f \in \mathbb{R}[x_1, \dots, x_s]$ of bounded degree, the test " $f(x) \in \Sigma^2$ " can be done in $\text{poly}(s)$ time.*

Let $\Sigma^{2,d}$ be those $f(x) \in \Sigma^2$ of degree at most d . Then $\Sigma^{2,d} \subseteq \mathcal{P}_+^d(\mathbb{R})$. To minimize $f(x)$ over \mathbb{R}^s , we find the largest $\lambda \in \mathbb{R}$ for which $f(x) - \lambda \in \Sigma^{2,d}$ via a binary search on λ and the proposition above. The value λ is a lower bound on $f(x)$ and in practice almost always agrees with the true minimum of f [25].

To minimize $f(x)$ over a set K constrained by polynomials, we need a few more tools. We could reduce the problem to minimizing a single polynomial, as mentioned in Section 6.1, but the following may work better in practice. We follow the presentation in [8].

Definition 6.5. The Algebraic Cone generated by elements $\beta_1, \dots, \beta_t \in \mathbb{R}[x_1, \dots, x_s]$ is the set,

$$\mathcal{A}(\beta_1, \dots, \beta_t) \stackrel{\text{def}}{=} \{f \in \mathbb{R}[x_1, \dots, x_s] \mid f = \eta + \sum_{I \subseteq [t]} \eta_I \prod_{i \in I} \beta_i\},$$

where η and the η_I are in Σ^2 , and $[t] = \{1, 2, \dots, t\}$.

Thus, the algebraic cone can be thought of as the set of all affine combinations of all possible products of polynomials β_1, \dots, β_t , where the coefficients of the affine combination are taken from Σ^2 .

Definition 6.6. The Multiplicative Monoid $\mathcal{M}(\beta_1, \dots, \beta_t)$ generated by $\beta_1, \dots, \beta_t \in \mathbb{R}[x_1, \dots, x_s]$ is the set of finite products of the β_i , including the empty product which we set to 1.

The key result is a simplified form of the Positivstellensatz [32]:

THEOREM 6.7. Given polynomials $\{f_1, \dots, f_{t_1}\}, \{g_1, \dots, g_{t_2}\}$ in $\mathbb{R}[x_1, \dots, x_s]$, the set

$$K \stackrel{\text{def}}{=} \{x \in \mathbb{R}^s : f_i(x) \geq 0, g_j(x) \neq 0, \forall i \in [t_1], j \in [t_2]\}$$

is empty iff $\exists F \in \mathcal{A}(f_1, \dots, f_{t_1})$ and $G \in \mathcal{M}(g_1, \dots, g_{t_2})$ for which $F + G^2$ is the zero polynomial.

Thus, for a set K described by f_i , and g_j of the form above, we consider $K' = K \cap \{x \in \mathbb{R}^s \mid \gamma - f(x) \geq 0, f(x) - \gamma \neq 0\}$. K' is empty iff $f(x) > \gamma$ for all $x \in K$.

Heuristics implemented in practice work by choosing a degree bound D , generating all $G \in \mathcal{M}(f - \gamma, g_1, \dots, g_{t_2})$ of degree at most D (there are at most t_2^D such G), and checking if there is an $F \in \mathcal{A}(\gamma - f, f_1, \dots, f_{t_1})$ for which $F + G^2 = 0$ via semidefinite programming. This is efficient for constant D , which usually suffices in practice. Better algorithms for special cases are based on alternative forms of the Positivstellensatz; see [27, 28].

7. CONCLUSION

We presented a novel approach to privacy where only gaining confidence in a sensitive fact is illegal, while losing confidence is allowed. We showed that this relaxation is significant and permits many more queries than with well-known approaches. In exchange, this gave us an opportunity to strengthen prior knowledge assumptions beyond current standards. Our hope is that work in this direction will help bridge the gap between theoretical soundness and practical usefulness of privacy frameworks.

One possible future goal is to obtain a better understanding of the families of sets and distributions that arise in practice, and to understand whether they admit efficient privacy tests. Another goal is to apply the new frameworks to online (proactive) auditing, which will require the modeling of a user's knowledge about the auditor's query-answering strategy.

Acknowledgements: We thank Kenneth Clarkson for bringing our attention to the Four Functions Theorem.

8. REFERENCES

- [1] R. Agrawal, R. J. Bayardo, C. Faloutsos, J. Kiernan, R. Rantzaou, and R. Srikant. Auditing compliance with a hippocratic database. In *Proc. VLDB*, pages 516–527, 2004.
- [2] R. Agrawal, J. Kiernan, R. Srikant, and Y. Xu. Hippocratic databases. In *Proc. VLDB*, pages 143–154, 2002.
- [3] R. Ahlswede and D. E. Daykin. An inequality for the weights of two families of sets, their unions and intersections. *Z. Wahrschein. und Verw. Gebiete*, 43:183–185, 1978.
- [4] S. Basu, R. Pollack, and M.-F. Roy. On the combinatorial and algebraic complexity of quantifier elimination. *J. ACM*, 43(6):1002–1045, 1996.
- [5] A. Blum, C. Dwork, F. McSherry, and K. Nissim. Practical privacy: The SuLQ framework. In *Proc. PODS*, pages 128–138, 2005.
- [6] B. Bollobás. *Combinatorics*. Cambridge Univ. Press, 1986.
- [7] J. Canny. Improved algorithms for sign determination and existential quantifier elimination. *Computer Journal*, 36(5):409–418, 1993.
- [8] C. Caramanis. Non-convex optimization via real algebraic geometry, 2001. http://web.mit.edu/~cmccaram/www/pubs/nonconvex_opt_review.pdf.
- [9] C. P. de Campos and F. G. Cozman. Computing lower and upper expectations under epistemic independence. In *Proc. 4th Intl. Symp. on Imprecise Probabilities and Their Apps.*, 2005.
- [10] I. Dinur and K. Nissim. Revealing information while preserving privacy. In *Proc. PODS*, pages 202–210, 2003.
- [11] C. Dwork and K. Nissim. Privacy-preserving datamining on vertically partitioned databases. In *Proc. CRYPTO*, pages 528–544, 2004.
- [12] A. Evfimievski, J. Gehrke, and R. Srikant. Limiting privacy breaches in privacy preserving data mining. In *Proc. PODS*, pages 211–222, 2003.
- [13] R. Fagin, J. Y. Halpern, Y. Moses, and M. Y. Vardi. *Reasoning About Knowledge*. The MIT Press, 1995. Paperback edition appeared in 2001.
- [14] R. Fagin, J. Y. Halpern, and M. Y. Vardi. A model-theoretic analysis of knowledge. *J. ACM*, 91(2):382–428, 1991.
- [15] S. Fujishige. *Submodular Functions and Optimization*, volume 58 of *Annals of Discrete Mathematics*. Elsevier, 2nd edition, 2005.
- [16] D. Grigoriev, E. de Klerk, and D. V. Pasechnik. Finding optimum subject to few quadratic constraints in polynomial time. In *Proc. Conf. on Effective Methods in Algebraic Geometry (MEGA)*, 2003.
- [17] J. Hintikka. *Knowledge and Belief*. Cornell University Press, 1962.
- [18] K. Kenthapadi, N. Mishra, and K. Nissim. Simulatable auditing. In *Proc. PODS*, pages 118–127, 2005.
- [19] S. Kripke. A semantical analysis of modal logic I: normal modal propositional calculi. *Zeitschrift für Mathematische Logik und Grundlagen der Mathematik*, 9:67–96, 1963. Announced in *J. of Symbolic Logic* 24, 1959, p. 323.
- [20] L. Lovász. Submodular functions and convexity. In A. Bachem, M. Grötschel, and B. Korte, editors, *Mathematical Programming – The State of the Art*, pages 235–257. Springer-Verlag, 1983.
- [21] G. Miklau and D. Suciu. A formal analysis of information disclosure in data exchange. In *Proc. SIGMOD*, pages 575–586, 2004.
- [22] R. Motwani, S. U. Nabar, and D. Thomas. Auditing SQL queries. In *Proc. ICDE*, 2008, to appear.
- [23] S. U. Nabar, B. Marthi, K. Kenthapadi, N. Mishra, and R. Motwani. Towards robustness in query auditing. In *Proc. VLDB*, pages 151–162, 2006.
- [24] P. A. Parrilo. Structured semidefinite programs and semialgebraic geometry methods in robustness and optimization, 2000. Ph.D. Thesis, California Institute of Technology.
- [25] P. A. Parrilo and B. Sturmfels. Minimizing polynomial functions. In *Algorithmic and Quantitative Aspects of Real Algebraic Geometry in Mathematics and Computer Science*, pages 83–100, 2001.
- [26] President's Information Technology Advisory Committee. Revolutionizing health care through information technology, 2004.
- [27] M. Putinar. Positive polynomials on compact semi-algebraic sets. *Indiana University Math Journal*, 42(3), 1993.
- [28] K. Schmüdgen. The k -moment problem for compact semialgebraic sets. *Annals of Math*, 289:203–206, 1991.
- [29] C. E. Shannon. Communication theory of secrecy systems. *Bell System Technical Journal*, 28-4:656–715, 1949.
- [30] N. Z. Shor. Class of global minimum bounds of polynomial functions. *Cybernetics*, 6:731–734, 1987.
- [31] N. Z. Shor and P. I. Stetsyuk. The use of a modification of the r -algorithm for finding the global minimum of polynomial functions. *Cybernetics and Systems Analysis*, 33:482–497, 1997.
- [32] G. Stengle. A Nullstellensatz and a Positivstellensatz in semialgebraic geometry. *Annals of Math*, 207:87–97, 1974.
- [33] G. H. v. Wright. *An Essay in Modal Logic*. North-Holland, 1951.