# INSIGHTS WITHOUT EXPOSURE

## Understanding 'differential privacy' in information security

**How should fraud examiners and legal professionals meet compliance standards but also keep individuals' and organizations' data private?**
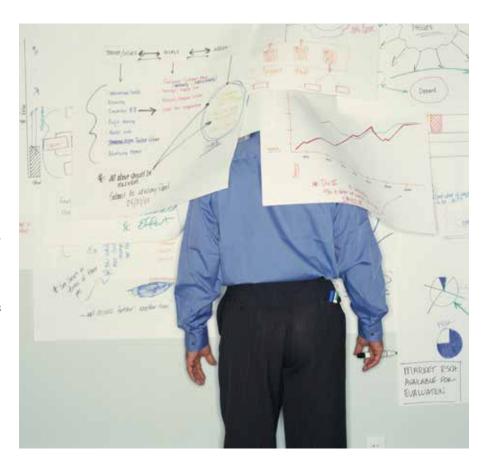"Differential privacy" is a new system of cybersecurity that proponents claim can protect data far better than traditional sanitizing or anonymizing methods.

Over the past two years as a columnist for *Fraud Magazine*, it's been a pleasure introducing new ideas, innovations and technologies to you — my colleagues. That's why I'm excited to devote this edition of "Innovation Update" to the concept of "differential privacy."

Differential privacy can securely limit algorithms so organizations can securely share private, sensitive data internally or among third parties. The concept isn't new. Mathematicians, cryptographers and academics have been discussing it for more than a decade. However, companies are now commercializing it for global fraud examinations and proactive compliance monitoring.

In September 2019, Google released the open-source version of the differential privacy library it uses in some of its products, such as Maps, according to Emil Protalinski, the author of "Google open-sources its differential privacy library," VentureBeat, Sept. 5, 2019, tinyurl.com/uhlautl.

"Differential privacy limits the algorithms used to publish aggregate information about a statistical database," Protalinski writes. "Whether you are a city planner, small business owner or software developer, chances are you want to gain insights from the data of



your citizens, customers or users. But you don't want to lose their trust in the process. Differentially private data analysis enables organizations to learn from the majority of their data without allowing any single individual's data to be distinguished or re-identified."

### Preventing reverse engineering of data

Here's a business case example to help clarify a challenge that differential privacy might meet. A company is managing a database containing sensitive personally identifiable information (PII),

such as customer credit cards, demographic and personal health information plus corporate product formulas and other forms of company intellectual property.

The company would like to release some statistics from this data to the public, a third-party vendor or joint-venture partner. However, the company has to ensure it's impossible for outsiders to reverse-engineer the released sensitive data. An outsider, in this example, would be an entity intending to reveal, or learn, at least some of the company's sensitive data elements.

Traditional approaches would most likely seek to simply anonymize the data (e.g., swap out customer names with random numbers) or even redact or delete sensitive fields in the data. However, if you have auxiliary information from other data sources coming into the repository, anonymization isn't sufficient because outsiders could reverse-engineer or cross-reference data sets to derive or recover the masked data.

For example, in 2007, Netflix released a dataset of its user ratings as part of a competition to see if anyone could outperform its collaborative filtering algorithm. The dataset didn't contain PII, but researchers were still able to breach privacy by cross-referencing other data sources to derive individual customer data. (See "Researchers Reverse Netflix's Anonymization," by Robert Lemos, SecurityFocus, Dec. 4, 2007, tinyurl.com/ux6k9sy.)

## Benefits of differential privacy

On my podcast, "The Walden Pond," I recently interviewed Ishaan Nerurkar, CEO of LeapYear Technologies Inc., a company that has applied differential

**COLUMNIST**
**VINCENT M. WALDEN, CFE, CPA**
MANAGING DIRECTOR, ALVAREZ & MARSAL'S DISPUTES AND INVESTIGATIONS PRACTICE

IN 2007, NETFLIX RELEASED A DATASET OF ITS USER RATINGS AS PART OF A COMPETITION TO SEE IF ANYONE COULD OUTPERFORM ITS COLLABORATIVE FILTERING ALGORITHM. THE DATASET DIDN'T CONTAIN PII, BUT RESEARCHERS WERE STILL ABLE TO BREACH PRIVACY BY CROSS-REFERENCING OTHER DATA SOURCES TO DERIVE INDIVIDUAL CUSTOMER DATA.

privacy research to develop a commercial platform for privacy-preserving computations on sensitive data. (See "Insights Without Exposure with Ishaan Nerurkar," The Walden Pond, at tinyurl.com/vo3nrrn.)

"Every regulated industry," Ishaan says, "whether it be in financial services, health care, telecom, aerospace and defense, government or industrial manufacturing — just to name a few — faces significant challenges using and sharing sensitive data."

"While there are some techniques for sanitizing data such as masking [or redacting] certain sensitive data fields, anonymizing the data or simply deleting key information, these techniques don't really lend [themselves] for today's level of analytics that require such valuable information in order to enhance predictive models or extract key insights required for effective decision-making," Ishaan says. "These old techniques either could reduce the value of the data, or worse — allow end users to perhaps even reverse-engineer the masked data and thus exposing the company to risk."

Ishaan further describes differential privacy as a technology that seeks to learn statistical patterns about the data without exposing underlying information. "It lets you run a statistic and build a model," Ishaan says. "But it won't allow for the exposure of a single underlying record that helped generate that model. … Think of it as a layer that sits on your database that allows the user to only gain access to the sensitive data through the differential privacy platform layers," he says. "Users are able to gain access to the database fields and select which fields need to be hidden in an easy-to-use interface."

Applications in which differential privacy algorithms can benefit an organization might include:

- Internal organization silos, in which various business units, such as human resources or legal, fear the compromise or theft of sensitive company data.

- Geographical borders, in which personal data privacy restrictions, such as the EU's General Data Protection Regulation (GDPR), might restrict data transferring between countries.

- Third parties or joint-venture partners, in which sales or customer

information might need to be exchanged but doing so could violate privacy policies or laws.

- Personal health care data, in which global patient data in clinical research needs to be analyzed to find a life-saving drug without violating individuals' specific medical information and data privacy details.

- Banking, in which financial institutions can build data platforms on traders that span across all their clients' information but don't expose any of that data related to a single client.

- Global investigations, litigation and compliance monitoring, in which information about a particular allegation or risk topic needs to be analyzed without violating an individual's data protection and privacy rights from one jurisdiction to the next.

"There have been significant advancements in differential privacy techniques" in global investigations and litigation "that can also apply to the identification of relevant information before a document production [in the context of a litigation]," according to the article, "Global Privacy Rules Intersect with Discovery Obligations," by Andy G. Gandhi, Mauricio Paez and Mark Kindy, *New York Law Journal,* Jan. 31, tinyurl.com/uz5ta8q.

## Open-source considerations

The technically and mathematically gifted have open-source options for using differential privacy. In addition to Google's version at the beginning of the column, academics frequently reference



APPLICATIONS WHERE DIFFERENTIAL PRIVACY ALGORITHMS CAN BENEFIT AN ORGANIZATION MIGHT INCLUDE PERSONAL HEALTH CARE DATA, IN WHICH GLOBAL PATIENT DATA IN CLINICAL RESEARCH NEEDS TO BE ANALYZED TO FIND A LIFE-SAVING DRUG WITHOUT VIOLATING INDIVIDUALS' SPECIFIC MEDICAL INFORMATION AND DATA PRIVACY DETAILS.

the development framework, "Pufferfish." See "Pufferfish: A Framework for Mathematical Privacy Definitions" by Daniel Kifer, of Penn State University. and Ashwin Machanavajhala, of Duke University, tinyurl.com/rd6j3rv, among several other public articles available on the internet.

According to the Duke University article, organizations can use the Pufferfish framework to create new privacy definitions that are customized to the needs of a given application. The goal of Pufferfish is to allow experts in a particular knowledge domain, who frequently don't have proficiency in privacy conventions, to develop rigorous privacy definitions for their data-sharing needs. Be forewarned: It's very technical reading.

## For fraud examiners

Is your curiosity piqued like mine? I hope so. In this small space, I can only give the rudiments of differential privacy. I encourage you to Google the topic and research further to see how organizations are applying it. I anticipate many more entities and sectors will be using differential privacy technologies as governments enact more data privacy regulatory laws, such as the GDPR and the California Consumer Privacy Act. ■ FM

---

**Vincent M. Walden, CFE, CPA**, is a managing director with Alvarez & Marsal's Disputes and Investigations Practice and is host of "The Walden Pond," a compliance podcast series. He welcomes your feedback. Contact him at vwalden@alvarezandmarsal.com.