

Efficient protocol to secure a user's private information when algorithms use it to recommend content

May 13 2022, by Adam Zewe



Researchers devise an efficient protocol to keep a user's private information secure when algorithms use it to recommend products, songs, or shows. Credit: Christine Daniloff, MIT

Algorithms recommend products while we shop online or suggest songs we might like as we listen to music on streaming apps.

These algorithms work by using [personal information](#) like our past purchases and browsing history to generate tailored recommendations. The sensitive nature of such data makes preserving privacy extremely important, but existing methods for solving this problem rely on heavy cryptographic tools requiring enormous amounts of computation and bandwidth.

MIT researchers may have a better solution. They developed a privacy-preserving [protocol](#) that is so efficient it can run on a smartphone over a very slow network. Their technique safeguards [personal data](#) while ensuring recommendation results are accurate.

In addition to user privacy, their protocol minimizes the unauthorized transfer of information from the database, known as leakage, even if a malicious agent tries to trick a database into revealing secret information.

The new protocol could be especially useful in situations where data leaks could violate [user privacy](#) laws, like when a [health care provider](#) uses a patient's medical history to search a database for other patients who had similar symptoms or when a company serves targeted advertisements to users under European privacy regulations.

"This is a really hard problem. We relied on a whole string of cryptographic and algorithmic tricks to arrive at our protocol," says Sacha Servan-Schreiber, a graduate student in the Computer Science and Artificial Intelligence Laboratory (CSAIL) and lead author of the paper that presents this new protocol.

Servan-Schreiber wrote the paper with fellow CSAIL graduate student

Simon Langowski and their advisor and senior author Srinivas Devadas, the Edwin Sibley Webster Professor of Electrical Engineering. The research will be presented at the IEEE Symposium on Security and Privacy.

The data next door

The technique at the heart of algorithmic recommendation engines is known as a nearest neighbor search, which involves finding the data point in a database that is closest to a query point. Data points that are mapped nearby share similar attributes and are called neighbors.

These searches involve a server that is linked with an online database which contains concise representations of data point attributes. In the case of a music streaming service, those attributes, known as feature vectors, could be the genre or popularity of different songs.

To find a song recommendation, the client (user) sends a query to the server that contains a certain feature vector, like a genre of music the user likes or a compressed history of their listening habits. The server then provides the ID of a feature vector in the database that is closest to the client's query, without revealing the actual vector. In the case of music streaming, that ID would likely be a song title. The client learns the recommended song title without learning the feature vector associated with it.

"The server has to be able to do this computation without seeing the numbers it is doing the computation on. It can't actually see the features, but still needs to give you the closest thing in the database," says Langowski.

To achieve this, the researchers created a protocol that relies on two separate servers that access the same database. Using two servers makes

the process more efficient and enables the use of a cryptographic technique known as private information retrieval. This technique allows a client to query a database without revealing what it is searching for, Servan-Schreiber explains.

Overcoming security challenges

But while private information retrieval is secure on the client side, it doesn't provide database privacy on its own. The database offers a set of candidate vectors—possible nearest neighbors—for the client, which are typically winnowed down later by the client using brute force. However, doing so can reveal a lot about the database to the client. The additional privacy challenge is to prevent the client from learning those extra vectors.

The researchers employed a tuning technique that eliminates many of the extra vectors in the first place, and then used a different trick, which they call oblivious masking, to hide any additional data points except for the actual nearest neighbor. This efficiently preserves database privacy, so the client won't learn anything about the feature vectors in the database.

Once they designed this protocol, they tested it with a nonprivate implementation on four real-world datasets to determine how to tune the algorithm to maximize accuracy. Then, they used their protocol to conduct private nearest neighbor search queries on those datasets.

Their technique requires a few seconds of server processing time per query and less than 10 megabytes of communication between the client and servers, even with databases that contained more than 10 million items. By contrast, other secure methods can require gigabytes of communication or hours of computation time. With each query, their method achieved greater than 95 percent accuracy (meaning that nearly

every time it found the actual approximate nearest neighbor to the query point).

The techniques they used to enable database privacy will thwart a malicious client even if it sends false queries to try and trick the server into leaking information.

"A malicious client won't learn much more information than an honest client following protocol. And it protects against malicious servers, too. If one deviates from protocol, you might not get the right result, but they will never learn what the client's query was," Langowski says.

In the future, the researchers plan to adjust the protocol so it can preserve privacy using only one server. This could enable it to be applied in more real-world situations, since it would not require the use of two noncolluding entities (which don't share information with each other) to manage the [database](#).

"Nearest neighbor search undergirds many critical machine-learning driven applications, from providing users with content recommendations to classifying medical conditions. However, it typically requires sharing a lot of data with a central system to aggregate and enable the search," says Bayan Bruss, head of applied machine-learning research at Capital One, who was not involved with this work. "This research provides a key step towards ensuring that the user receives the benefits from nearest neighbor search while having confidence that the central system will not use their data for other purposes."

More information: Private Approximate Nearest Neighbor Search with Sublinear Communication. eprint.iacr.org/2021/1157.pdf

This story is republished courtesy of MIT News

(web.mit.edu/newsoffice/), a popular site that covers news about MIT research, innovation and teaching.

Provided by Massachusetts Institute of Technology

Citation: Efficient protocol to secure a user's private information when algorithms use it to recommend content (2022, May 13) retrieved 30 November 2023 from <https://techxplore.com/news/2022-05-efficient-protocol-user-private-algorithms.html>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.