24-09-2024

How to study data that you don't have access to?

Irina Bejan (OpenMined)





We are a non-profit foundation who builds open source technology to help ease the challenges of researching sensitive data.

OpenMined

KEY PROBLEM:

Our ability to answer important questions is limited

because we can't access existing data



KEY PROBLEM:

Our ability to answer important questions is limited

because we can't access existing data

Another country

Another org

Another department



Use Case 1: Global Trade



Canada's exports to the USA **usually** ~match USA's imports from Canada



... but sometimes goods *disappear*



... other times goods *multiply*



key premise: us can't access most of the world's data about its imports... ...because it's stored in other country



CA can't access most of the world's data about its imports... ...because it's stored in other country



Economic Surveillance

Trade imbalances

Market inefficiencies

Price fluctuations

Due of the lack of access, important questions are left unanswered

Supply shortages

Foreign Policy misalignment

Tarrifs and taxes





Service Canada

Canada



*







WHY CAN'T WE TRACK GLOBAL TRADE TOGETHER?

The world's national trade statistics orgs can't put all their data into one place





Use Case 2: Money Laundering



UNITED NATIONS OFFICE on DRUGS & CRIME:

\$1.6 trillion in criminal proceeds are laundered annually



WHY CAN'T WE STOP IT? Citi can't access most of the world's data about where its money comes from... ...because it's stored in other banks







Scotiabank[®]



SOCIETE

GROUPE BPCE

ussia



Crédit 🔂 Mutuel

WHY CAN'T WE TRACK MONEY LAUNDERING TOGETHER?

CIERCOR SUNTRUST SUNTRUST SUNTRUST LEMorgan J.P.Morgan The world's banks can't put all their ledger data into one place to follow the money

MIZUHO



SMFG



MUFG

Use Case 3: Emergency Response



UNITED NATIONs:

In 2020, approximately 270 million people were affected by disasters worldwide.



The IFRC states that 134 million people annually need humanitarian assistance due to natural disasters





Level(3)



Sprint



AT&T

Axtel 🔅

WHY DON'T WE KNOW WHERE TO SEND HUMANITARIAN AID?



TELECOM





comcast

···**T**··Mobile·

Time Warner

We could continue with more ...

Use Case 1: Global Trade Use Case 2: Money Laundering Use Case 3: Emergency Response Use Case 4: Breast Cancer Use Case 5: Algorithmic Accountability Use Case 6: Disinformation





Because..

...data access is *really hard*







How are we solving it?





Privacy Enhancing Technologies (PETs)













PETS MAKE IT POSSIBLE TO:

answer a question using data owned by someone else



In Another Org

In Another Dept.









...but we don't yet have a <u>car</u>.



KEY QUESTION:

What is the "car" of PETs?



OpenMined / PySyft	♥ Sponsor ⑦ Used by - 16 Unwatch - 167 ★ Unstar 3,893 Y For	k 869
<> Code (1) Issues 162 (1) P	Pull requests 17 🔹 Actions 🕕 Security 🔟 Insights 🔅 Settings	
library for encrypted, privacy p	preserving deep learning	Edit
deep-learning secure-computation	n pytorch privacy cryptography Manage topics	
© 5,337 commits	branches 🛇 5 releases 🛷 1 environment 👫 180 contributors 🎄 Apache	-2.0
Branch: dev - New pull request	Create new file Upload files Find File Clone or do	wnload 🗸
ivmaneuro and robert-wagner r	m peopless integration test (#2620)	
		Jays ago
.github/ISSUE_TEMPLATE	Update issue templates 4 mor	iths ago
art	Improve the diagram 3 mor	iths ago
docker-image	delete unecessary package installation, numpy comes already with pysyft 4 more	ths ago
docs	bumpversion 0.1.27a1 -> 0.1.28a1 (#2619) 13 c	lays ago
examples	Improvements to the Federated Recurrent Neural Network notebook (#2613) 18 c	lays ago
images	Add files via upload 3 mor	ths ago
syft	Implementing Protocol (#2605) 7 c	lays ago
test	rm needless integration test (#2629) 3 c	lays ago
.flake8	changed ignore to select 6 mor	oths ago

PySyft 1

The "car" of PETs




























Download Answers















3 D

Download Answers





Notice what's missing!

- **Privacy Risk:** researcher only sees results
- Security Risk: researcher doesn't have internal access to secure systems
- IP Risk: researcher only sees results (no other IP), as checked by admin.











UN PET Lab



PETs for NSOs

- Mandate to serve the public interest by
 - Protecting citizen data
 - Maximising the public benefits of the data
- To maximise the public benefits, NSOs need third-party data, such as:
 - o administrative data from other agencies
 - o data held in the private sector







PETS MAKE IT POSSIBLE TO:

answer a question using data owned by someone else



In Another Org

In Another Dept.



But there's another problem...



PROBLEM

To jointly study private data owned by two or more organisations, the data needs to be put on <u>the same</u> <u>computer.</u>



PROBLEM

They need **mutual secrecy**



StatCan & US Census launch datasites..





























Solution: Launch a secure enclave



How does it work?



Step 1: External Researcher proposes a project



Step 2: Admins review the code



Step 3: Admins approve the code



 \diamondsuit

Step 4: Enclave is verified



 \diamondsuit

Step 4: Enclave is verified



Step 5: Execution happens



Step 5: Execution happens



Step 6: Researcher gets answers


Ongoing Project 3: UN Global Platform to enable cross-border joins on global trade







Cont

enable

cros

A New Model for International, Privacy-Preserving Data Science

Monday, June 03, 2024 - 4:05 pm-4:20 pm Curtis Mitchell, xD, US Census Bureau

Abstract:

Currently when data analysis is performed between National Statistical Organizations (NSOs) such as the US Census Bureau and Statistics Canada, a complex series of arrangements must be agreed to that creates severe vet important restrictions on how and by whom the required data is accessed, thus increasing burden and time.

Here we demonstrate a new approach using remote, privacy-preserving processes via a collaboration between multiple NSOs in conjunction with the United Nations Privacy-Enhancing Technologies Lab (UN PET Lab). The proof-of-concept involves using the open-source data science platform PySyft and establishing the cloud infrastructure necessary such that nodes hosted by the US Census Bureau and other NSOs are facilitated by a network gateway hosted by the UN PET Lab. This architecture enables a private join on synthetic data representing realistic trade data from UN Comtrade, without each NSO needing to directly access the other NSO's data. It also enables investigations into key policy and governance questions as these technologies mature.

We believe this project will be an important milestone towards enabling privacy-preserving and remote data science between international government entities and uncovering future aspects of privacy policy and governance.



Curtis Mitchell is an Emerging Technology Fellow on the xD team at the US Census Bureau where he is contributing to a variety of projects involving privacy-enhancing technologies, responsible artificial intelligence, and modern web applications. He has over 15 years of experience in software- and data-related roles at small startups, large corporations, and open-source communities. Prior to joining the Census Bureau, he worked at NASA's Ames Research Center.

Connect: LinkedIn

BibTeX



Ongoing Project 4: enable joins between public & sector parties to study safety of frontier AI models



Goal: researchers can study any non-public data to answer important questions





Thank you for your attention!

Check out our work at openmined.org

