# Federated Machine Learning : Concept and Applications

**HUANG liu**

e0575772@u.nus.edu

Background

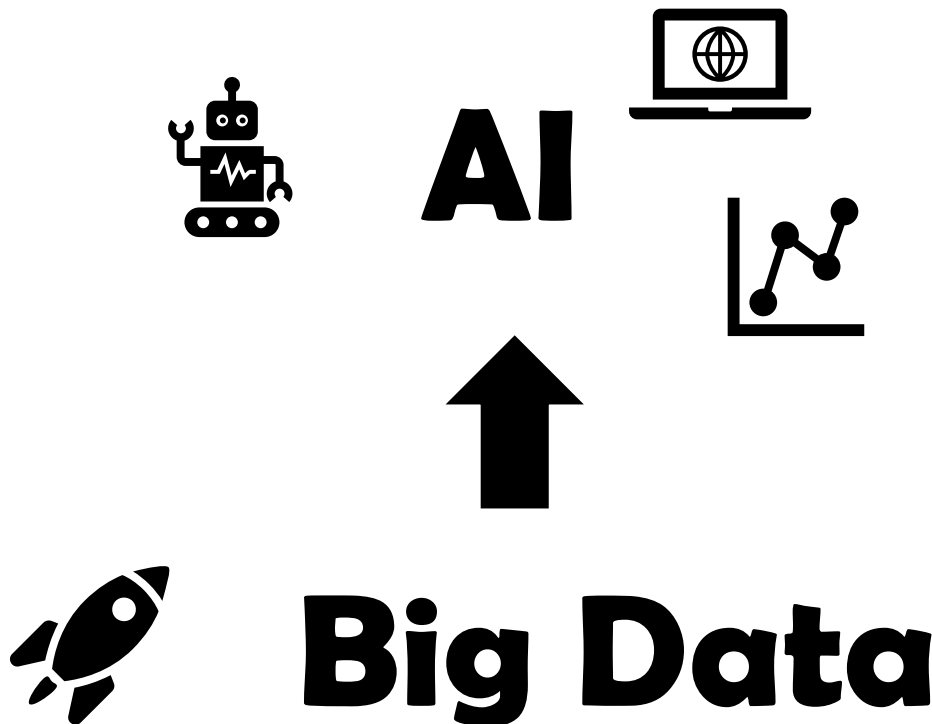Federated Learning

Related Work

Application

Conclusion

# 1 Background

# Growing Interest in AI

- Alpha GO
- Training dataset: 300,000 games

# Real World Situations

Limited amount of data
Low quality

Fuse data

High cost to exchange:
Strict laws
Complicated administrative procedures

Legally

Industry competition
Privacy security

Ensure security

# Major challenges

## Solution

| Data Island |

| Data Security |

→ Federated Learning

# 2 Federated Learning

# Definition

Federated Learning is a machine learning setting where the goal is to train a high-quality centralized model while training data remains distributed over a large number of clients each with unreliable and relatively slow network connections. [1]

[1] Federated Learning: Strategies for Improving Communication Efficiency. CoRRabs/1610.05492(2016). arXiv:1610.05492 http://arxiv.org/abs/1610.05492

# Privacy of Federated Learning

- Secure Multi-party Computation (SMC).

- Differential Privacy

- Homomorphic Encryption

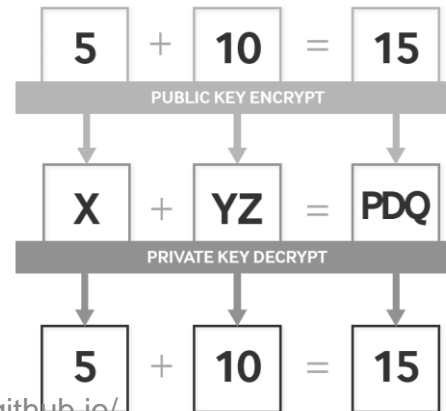Differentially Private Federated
Learning: A Client Level Perspective

$$w_{t+1} = w_t + \frac{1}{m_t} \left( \overbrace{\sum_{k=0}^{m_t} \triangle w^k / \max(1, \frac{\|\triangle w^k\|_2}{S})}^{\text{Sum of updates clipped at } S} + \overbrace{\mathcal{N}(0, \sigma^2 S^2)}^{\text{Noise scaled to } S} \right)$$

Gaussian mechanism approximating sum of updates



Figure 4: Detailed description of the Secure Aggregation protocol. Red, underlined parts are required to guarantee security in the active-adversary model (and not necessary in the honest-but-curious one).

SMPAI: Secure Multi-
Party Computation for
Federated Learning

https://images.a
pp.goo.gl/pCjC7
bf5FzRKaDsz8

https://huangliu0909.github.io/                9

# Categorization

- Horizontal Federated Learning



Similar Features
Various Samples

Eg: two banks in
different cities

# Categorization

- Horizontal Federated Learning
- Vertical Federated Learning



Similar Samples
Various Features

Eg: bank and
shopping mall in
the same city

# Categorization

- Horizontal Federated Learning
- Vertical Federated Learning
- Federated Transfer Learning



Various Samples
Various Features

Eg: bank in China
and shopping mall in
Singapore

"Similarity"

# Architecture of HFL



No information leakage between any parties.
Independent from specific ML algorithms
All participants will share the final model parameters.

# Architecture of VFL

Suppose A has training data Xa while B has training data Xb and labels y. We want to model how Xa and Xb jointly influence the value of label.

Since A and B cannot exchange data directly, we need a third party, C to help with the model training.



a

b

# Architecture of FTL

- The same architecture as VFL
- Differ in detail when trying to find the common representation among the parties
- Incentive mechanism: after the model is built, the local model's performance depends on how much this party contribute to the whole federated system.

# **3** **Related Work**

# Privacy-preserving ML

Federated learning can be considered as privacy-preserving decentralized collaborative machine learning.
Most of the privacy protection techniques using in privacy-preserving machine learning can be applied in Federated learning



$\triangle W = Aggr ( \triangle W1^{dp} + \triangle W2^{dp} + ... + \triangle Wn-1^{dp} + \triangle Wn^{dp} )$

Untrusted Server

$\triangle W1^{dp}$   $\triangle W2^{dp}$   $\triangle Wn-1^{dp}$   $\triangle Wn^{dp}$

$\triangle W$   $\triangle W$   $\triangle W$   $\triangle W$

.........

Hospitals

https://huangliu0909.github.io/

# Distributed Machine Learning



Parameter Server: $W_{i+1} = \frac{1}{4} \sum_{j=1}^{4} W_{i+1,j}$

$W_{i+1,1}$  $W_{i+1,2}$  $W_{i+1,3}$  $W_{i+1,4}$

$W_i$  $W_i$  $W_i$  $W_i$

Machine 1   Machine 2   Machine 3   Machine 4

https://code-it.ro/an-introductory-guide-on-distributed-training-of-neural-networks/

Non-IID local data

# Distributed Machine Learning

- DML: the parameter server allocates data on distributed working nodes and computes model parameters in a scheduled way.

- FL: each working node i.e. data holder, can independently decide when and how to join federated learning.

  Focus on privacy protection

# Edge Computing



CLOUD

EDGE
Service delivery
Computing offload
IoT management
Storage & caching

Edge Node

Edge Node

https://images.app.goo.gl/1fNMc1QT1Lw8jpcQ7

Federated learning can provide protocols of implementation details for edge computing, can work as an operation system for edge computing.

# Federated Database Systems

- Systems that integrate multiple database units and manage the integrated system as a whole

# Federated Database Systems

- VS Distributed database system

The data in each database unit is heterogeneous.

- VS Federated learning

Similar in terms of type and storage in data.

Focus on basic operations of data rather than training a machine learning model.

No privacy protection.

**4** Application

# Smart Retail

Goal: provide customers with personized services, such as product recommendation and sales service

| purchasing power | user's preference | information of product |
|:---:|:---:|:---:|
|  |  |  |
| Bank Saving | Social Media | e-shops |

# Smart Retail

Problem: data are scattered and heterogeneous

Solution: federated learning & transfer learning

"Heterogeneous data are any data with high variability of data types and formats. They are possibly ambiguous and low quality due to missing values, high data redundancy, and untruthfulness. It is difficult to integrate heterogeneous data to meet the business information demands."
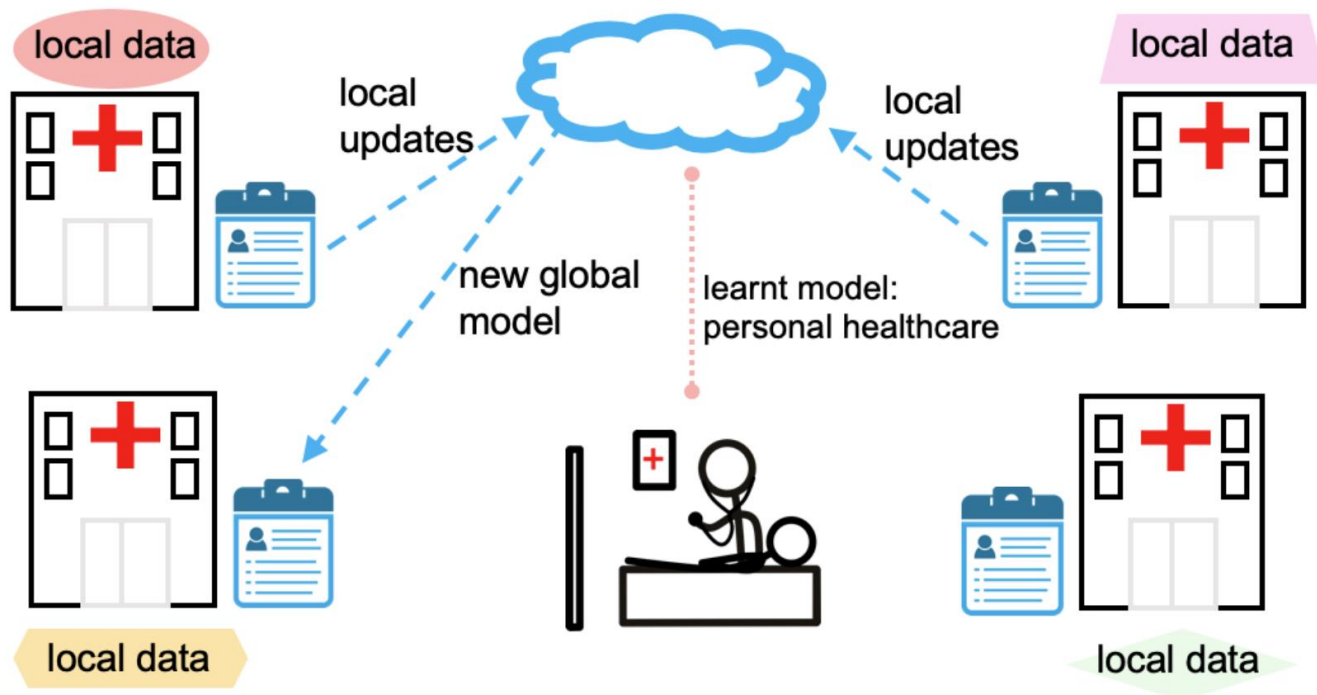
Wang, L. (2017). Heterogeneous Data and Big Data Analytics. *Automatic Control and Information Sciences*, *3*(1), 8-15.

✓ cross-enterprise

✓ cross-data

✓ cross-domain

# Smart Healthcare

Problem: data island & data security

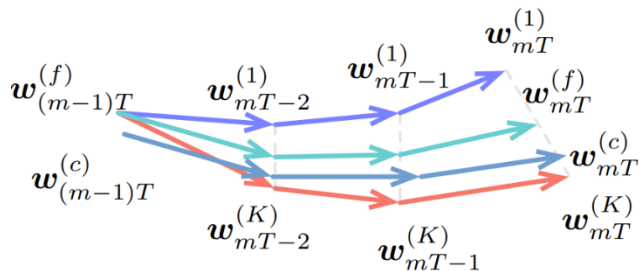Solution: horizontal federated learning

**5** **Conclusion**

"It is expected that in the near future, federated learning would break the barriers between industries and establish a community where data and knowledge could be shared together with safety, and the benefits would be fairly distributed according to the contribution of each participant."
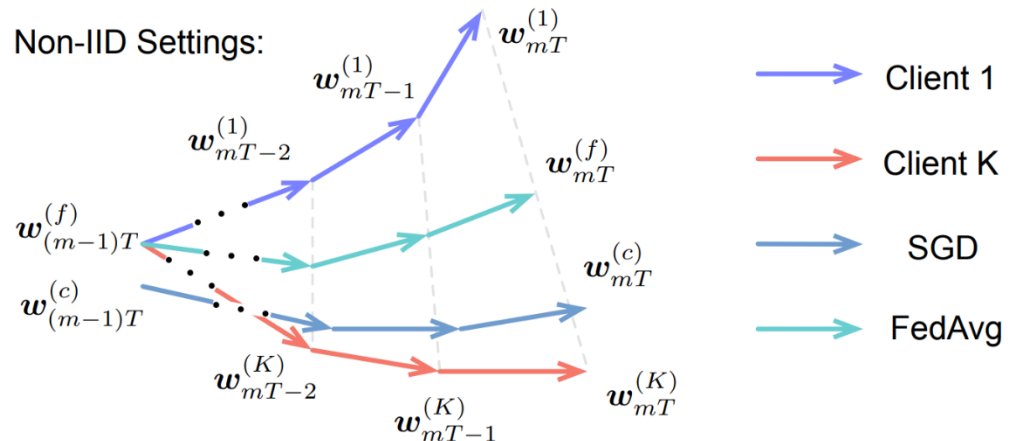
# My insights

- Non-IID data: adaptive optimizer

In the computing of each step, the learning rate is modified according to the historical gradient. The main advantage is to individually learn from the local distribution.



YueZhao,MengLi,LiangzhenLai,NaveenSuda,DamonCivin,andVikasChandra.2018. Federated Learning with Non-IID Data

# My insights

- Asynchronous Federated Learning

All the framework above assume little delay for each nodes' model transferring to the server, i.e. the server has to collect all gradients before updating global model.

What if server wait for too much time for collecting all gradients in one round?

Every time the server receives a model from A, the server updates the global model and immediately sends the new model back to A. So there're less time waste in one round and more efficiency.

# THANKS