# Federated Learning Algorithms

Andrew Zhang, Jiayi Tong, Yong Chen

University of Pennsylvania Biostatistics Dept.

## Introduction

The objective of this research is to understand and experiment with different approaches to federated learning. Federated Learning algorithms are unique in that they are able to train machine learning models while keeping user-level data secure and private. Normally, Machine learning models must centralize training data to create the model. However with federated learning, the data never needs to be moved, instead only summary statistics are transmitted. The use cases of this range from using patient data in hospitals to customer data in banks.

## Approaches

As a part of this project, I had experimented with various approaches to federated machine learning

1. PennCil Group's federated Logistic regression algorithm. This algorithm used taylor approximations and hessians to create a loss function that depends only on summary level statistics.

2. Distributed SVM techniques are primarily for high dimensional datasets. However in practice and in industry, SVM's have fallen out of favor due to inefficiency in training time.

3. XGBoost is a novel decision tree learning algorithm which is both quick and accurate, leading to wide adoption in industry.

## Ensemble Learning

Perhaps the simplest federated learning implementation is ensemble learning where a model is trained at each site and then voting occurs between the models. This technique is easy to setup and trains quickly but runs under the assumption that all data sites have a sufficient amount of data to train an efficient model. It also runs under the assumption that each location has data of the same distribution.

## Implementation of Ensemble XGBoost



XGBoost on one site's data (61% accuracy)

2 sites data regular xgboost (71% accuracy)



Federated: Voting between two sites (70% accuracy)

Federated Voting between 2 sites (logistic reg voting) The log reg part trained with central site data

When testing it on the central site, the results are very good

But when it is tested on test set, we see that it is overfitting and actually hurting overall effectiveness

## Summary

This Summer yielded many fruitful projects and learning experiences.

1. I completed the Stanford Online Machine Learning certificate.

2. I became familiar with popular tools such as Tensorflow, Sklearn, XGBoost, LightBGM, and others.

3. I conducted literature review of distributed SVM methods and investigated how it stacks up with alternative federated learning approaches.

4. I pivoted to XGBoost and initiated research into both the open-source code and Berkeley's Rise Lab's research into creating histogram based, Federated XGBoost.

## Web Application Project

Alongside Machine Learning projects, I helped develop a data sharing web application with Flask and MongoDB. This website aims to facilitate data transfer between sites for users of PDA (privacy preserving distributed algorithms) by allowing for summary statistics to be shared in json format.

I worked on the REST API, database, and parts of the frontend.