

# Vertical Federated Learning with Gossip Protocol

DVPL-Katyusha with FastMix

---

Ivan Toropin, Pyotr Lisov, Aleksandr Besnosikov

21.05.2024

MIPT, Russia

# Table of contents

1. Introduction
2. Motivation
3. Related Works
4. The Problem Statement
5. Results
6. Conclusion
7. Literature

# Introduction

The problem of time in such fields as ML and DL is more relevant than ever. Today it is common that teaching large neural networks requires several days or even weeks. **Federated Learning** is a new solution for this issue. It uses several machines, (usually called servers, agents, nodes, edges, etc.), that can compute the algorithm separately and then exchange data with each other.

# Introduction

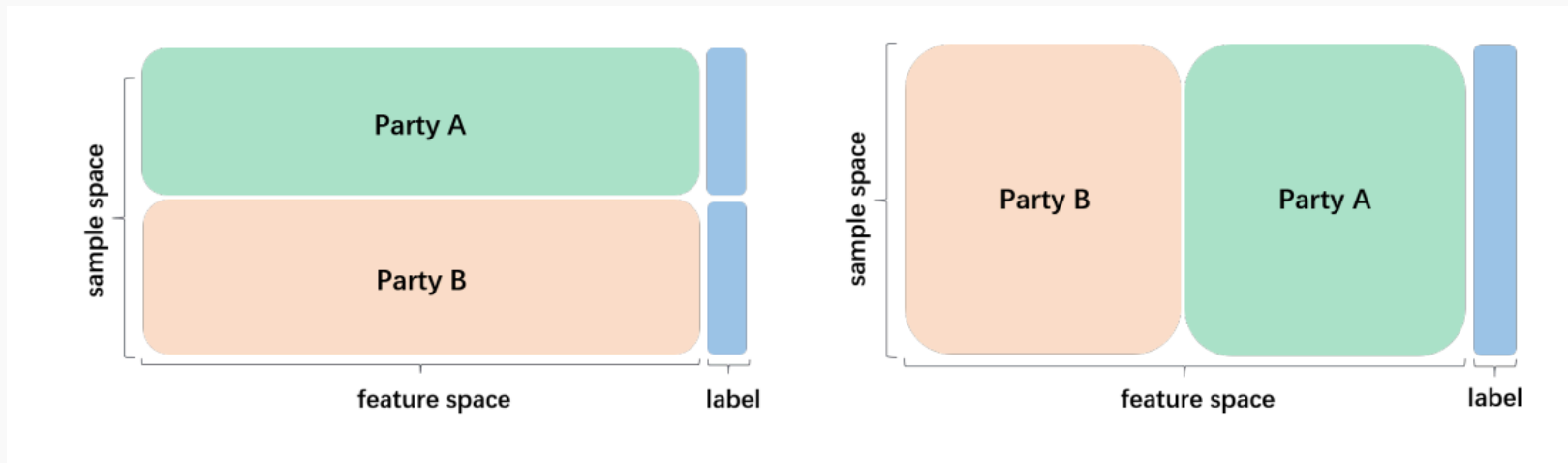


Figure 1: HFL vs VFL

Federated Learning can be further divided into two categories: **horizontal**, i.e each node share the same feature space while holding different samples, and **vertical**, i.e each node share the same samples/users while holding different features

# Motivation

The main reason to use VFL is that often data owned by different parties have the same sample IDs but disjoint subsets of features. Such scenario is common in the industry applications of collaborative learning, such as medical study, financial risk, and targeted marketing. For example, E-commerce companies owning the online shopping information could collaboratively train joint-models with banks and digital finance companies that own other information of the same people such as the average monthly deposit and online consumption, respectively, to achieve a precise customer profiling.

# Motivation

It is also quite common that different parties can't communicate with each other directly (this case is called fully connected networks). Instead, they have a small range of their 'neighbors'. Thus, another important problem is the organization of communication between different nodes in sparse networks. One of the ways to model such communications is the **gossip protocol** or epidemic protocol, which we decided to use. It is a procedure or process of computer peer-to-peer communication that is based on the way epidemics spread.

## Related Works

We decided to use the **DVPL-Katyusha** (Sergey Stanko & Timur Karimullin, 2024) algorithm as a basis. It is a VFL implementation of **L-Katyusha** (Kovalev et al., 2020), developed for fully connected networks. In order to modify it for sparse ones we used a variant of gossip protocol called **FastMix** (Haishan Ye et al., 2023) in stead of **AllReduce** (Ernie Chan et al., 2007)

# The problem statement

Like in the original paper (Sergey Stanko & Timur Karimullin, 2024) we consider the problem of **minimization of the linear loss function**, distributed among  $n$  devices. This can be mathematically written as:

$$\min_{x \in \mathbb{R}^d} \left[ f(x) := \frac{1}{s} \left\| \sum_{i=1}^n A_i x_i - b \right\|^2 \right]$$

We assume that  $f(x)$  is  $\mu$ -**strongly convex** and  $L$ -**smooth**.

Our goal is to prove the convergence of the new algorithm, show that it is at least as quick as the original one.



# Results

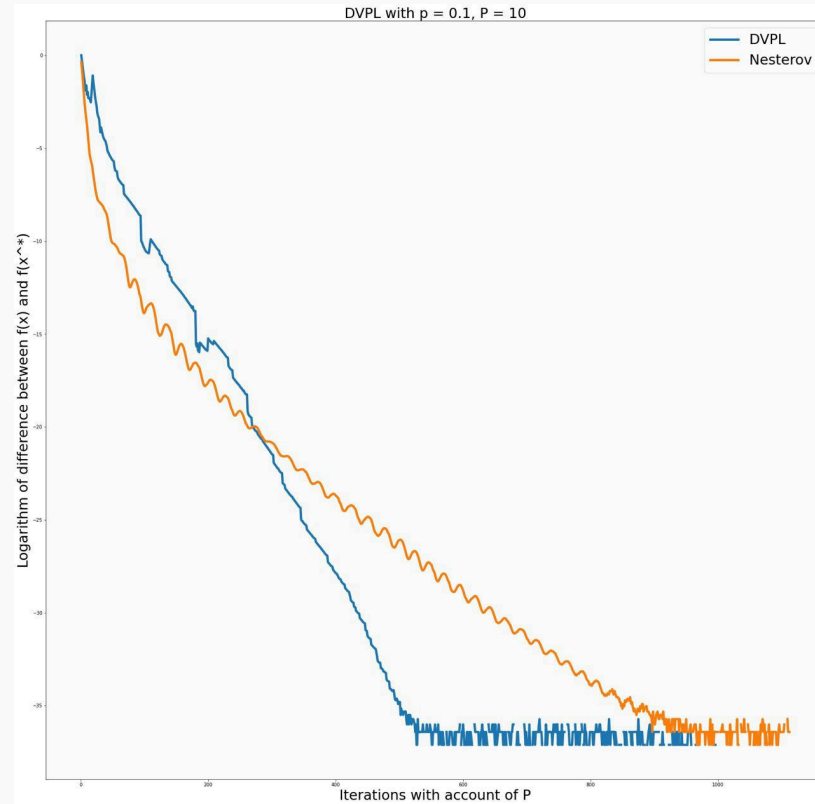


Figure 2: HFL vs VFL

# Results

Before attempting to prove theoretical convergence, a numerical simulation was conducted. Thus, we estimated the optimal parameters for the convergence of the method.

# Results

Finally, choosing right parameters we can get an estimation of the number of iterations.

Let  $p = \frac{b}{s}$ ,  $\theta_1 = \min\left\{\sqrt{\frac{2\sigma s}{3b}}, \frac{1}{2}\right\}$ ,  $\theta_2 = \frac{1}{2}$ , then it will require

$$K = \mathcal{O}\left(\frac{s}{b} \sqrt{\frac{L}{\mu} \left(\frac{1}{s} + \frac{\sum_{j=1}^s L_j^2}{\mu^2}\right)} \log \frac{1}{\varepsilon}\right) \sim \mathcal{O}\left(\sqrt{\frac{s}{b}} \sqrt{\frac{L^3}{\mu^3}} \log \frac{1}{\varepsilon}\right)$$

iterations to achieve  $\mathbb{E}[\Psi^{k+1}] < \varepsilon \Psi^0$ , where  $\Psi^k = Z^k + Y^k + W^k$  is the **Lyapunov function**. This concludes our proof.

Our main result:

$$K = \mathcal{O} \left( \sqrt{\frac{s}{b}} \sqrt{\frac{L^3}{\mu^3}} \log \frac{1}{\varepsilon} \right)$$

Compared to the original paper:

$$K = \mathcal{O} \left( \frac{\sqrt{s}}{b} \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} \right)$$

- [1] Sergey Stanko & Timur Karimullin, Accelerated Methods with Compression for Horizontal and Vertical Federated Learning, 2024
- [2] Kovalev, S. Horvath, and P. Richtarik. Don't jump through hoops and remove those loops: Svrg and katyusha are better without the outer loop, 2020
- [3] Haishan Ye, Luo Luo, Ziang Zhou, and Tong Zhang. Multi-consensus decentralized accelerated gradient descent. arXiv preprint arXiv:2005.00797, 2020
- [4] E. Chan, M. Heimlich, A. Purkayastha, R. V. D. Geijn. Collective communication: theory, practice, and experience, 2007