

Accelerated Methods with Compression for Horizontal and Vertical Federated Learning

Sergey Stanko, Timur Karimullin, Aleksandr Beznosikov

Moscow Institute of physics and technology

Introduction

To address the problem of the time-consuming process of function minimization, the community came to distributed algorithms, where the calculation process is divided among different devices. Such parallel computation can be used in situations, where data is distributed across several machines, as in the case of federated learning approach. The latter can be divided into two different regimes. One of these is horizontal federated learning, where each worker possesses their own collection of samples, but share the same set of features. A different way of data division is considered in the vertical case, where every device has a unique set of features of the same samples.

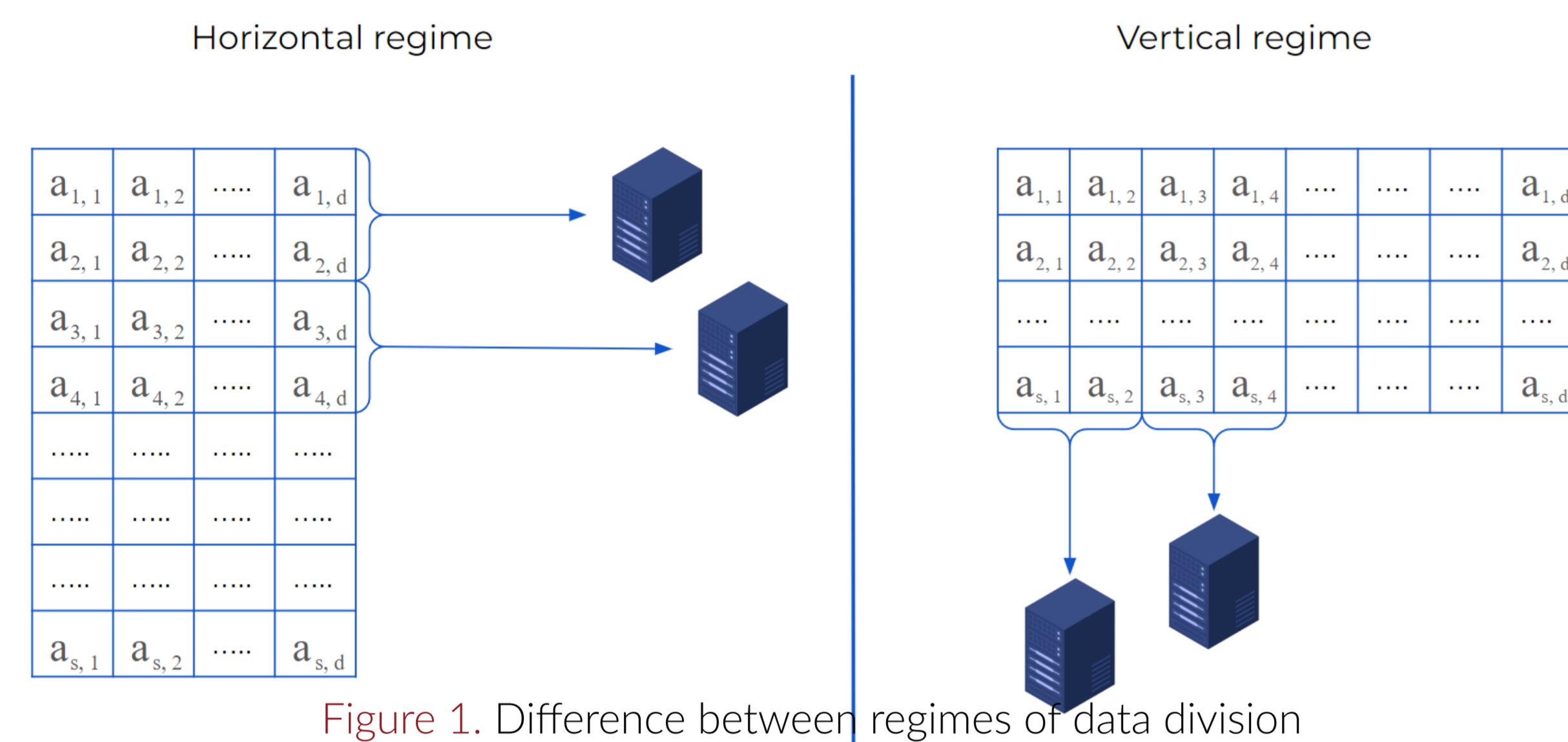


Figure 1. Difference between regimes of data division

Compression

Despite the considerable practical utility of distributed methods, they are not without significant drawbacks. In particular, the parallelisation of a task does not necessarily result in an optimal reduction in time. That means that having n devices does not accelerate task by n times. This happens because of limited ability of networks to exchange information. Thus, the key bottleneck of parallel computation is the communication part. There have been considered several ways of dealing with this issue [2], but in our paper we concentrate solely on reducing communication cost of each iteration by decreasing the size of sending information also known as compression technique [1].

The main results

Table 1: Summary of bounds for iteration complexities for finding an ε -solution. Convergence is measured by the distance to the solution.

Regime	Reference	Iteration complexity
Horizontal	QSGD Alistarh et al. [2017] ⁽¹⁾	$\mathcal{O}\left(\frac{L}{\mu} \left(1 + \frac{\omega}{n}\right) \log \frac{1}{\varepsilon} + \frac{\sigma_*^2 \omega}{\mu^2 \varepsilon}\right)$
	DIANA Mishchenko et al. [2023]	$\mathcal{O}\left(\left(\frac{L}{\mu} \left(1 + \frac{2\omega}{n}\right) + \omega\right) \log \frac{1}{\varepsilon}\right)$
	VR-DIANA Horváth et al. [2019]	$\mathcal{O}\left(\left(\frac{L}{\mu} \left(1 + \frac{\omega}{n}\right) + \omega\right) \log \frac{1}{\varepsilon}\right)$
	ADIANA Li et al. [2020]	$\mathcal{O}\left(\left(\left(\frac{L}{\mu} \left(1 + \sqrt{\frac{\omega}{n} + \frac{\omega}{n}}\right) \omega\right) + \omega\right) \log \frac{1}{\varepsilon}\right)$
	ADIANA He et al. [2024]	$\mathcal{O}\left(\left(\frac{L}{\mu} \left(\sqrt{\frac{\omega^2}{n} + 1}\right) + \omega\right) \log \frac{1}{\varepsilon}\right)$
	MARINA Gorbunov et al. [2022] ⁽²⁾	$\mathcal{O}\left(\left(\frac{L}{\mu} \left(1 + \frac{\omega}{\sqrt{n}}\right) + \omega\right) \log \frac{1}{\varepsilon}\right)$
	MASHA Beznosikov et al. [2023] ⁽³⁾	$\mathcal{O}\left(\left(\frac{L}{\mu} \sqrt{\left(w + \frac{\omega^2}{n}\right) + \omega}\right) \log \frac{1}{\varepsilon}\right)$
	Three Pillars Algorithm Beznosikov and Gasnikov [2023] ^(3, 4)	$\mathcal{O}\left(\left(\frac{L}{\mu} \sqrt{n} + n\right) \log \frac{1}{\varepsilon}\right)$
	Algorithm 1(This paper)	$\mathcal{O}\left(\left(\frac{L}{\mu} \left(\sqrt{\frac{\omega^2}{n} + 1}\right) + \omega\right) \log \frac{1}{\varepsilon}\right)$
Vertical	AVFL Cai et al. [2022]; CE-VFL Sun et al. [2023]; SecureBoost+ Chen et al. [2021]; eHE-SecureBoost Xu et al. [2021];	No theoretical results
	CVFL Castiglia et al. [2023]	No concrete number of iterations ⁽⁵⁾
	Algorithm 2(This paper) ⁽⁴⁾	$\mathcal{O}\left(\left(\sqrt{\frac{L\omega}{\mu}} + s\right) \log \frac{1}{\varepsilon}\right)$

⁽¹⁾ with tuned step size, no convergence with fixed step size; ⁽²⁾ under Polyak-Lojasiewicz condition; ⁽³⁾ for variational inequalities and saddle point problems; ⁽⁴⁾ for special compressor only; ⁽⁵⁾ for special compressor convergence rate is $\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$, no guarantees in the case of arbitrary compressor;

Notation: μ = constant of strong convexity, L = smoothness constant, ω = compression constant (see Definition 1.1), n = number of workers, σ_* = sum of workers' gradients in the optimal point, s = total number of samples.

Problem statement

We state the standard distributed learning problem, which, formally, can be written in the following form.

$$\min_{x \in \mathbb{R}^d} \left[f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) \right].$$

We assume functions to have the following properties.

Definition 2: The functions $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ are L -smooth for some $L > 0, \forall i \in \overline{1, n}$:

$$f_i(y) \leq f_i(x) + \langle \nabla f_i(x), y - x \rangle + \frac{L}{2} \|y - x\|^2, \forall x, y \in \mathbb{R}^d.$$

Definition 3: The function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is μ -strongly convex for some $\mu > 0$:

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2, \forall x, y \in \mathbb{R}^d.$$

Compressor properties

Compression mathematically can be represented in the form of a vector function. In our paper, we consider them having the following properties:

Definition 1 We say that the compression operator Q is unbiased, if

$$\mathbb{E}[Q(x)] = x, \forall x \in \mathbb{R}^d.$$

We also assume that there exists a constant $\omega \geq 0$, such:

$$\mathbb{E}[\|Q(x) - x\|^2] \leq \omega \|x\|^2, \forall x \in \mathbb{R}^d.$$

DHPL-Katyusha

The DHPL-Katyusha, the algorithm for the horizontal data division regime, is inspired by the original L-Katyusha [3]. We use the variance reduction technique for compressor's stochastics, thus in our method, every device calculates the compressed gradient difference and broadcasts it to other workers. After that, the new points are found using the momentum part, similar to L-Katyusha.

Algorithm 1 DHPL-Katyusha

Parameters: $\theta_1, \theta_2, \gamma \in (0, 1]$, probability $p \in (0, 1]$ (every worker has the same random seed for random, connected with p).

Initialization: Choose $y^0 = w^0 = z^0 \in \mathbb{R}^d$, stepsize $\eta = \frac{\theta_2 \gamma}{(1+\theta_2)\theta_1}$, set $\tilde{L} = L \max\{\frac{\omega}{n}, 1\}$, $\sigma = \frac{\sigma_*}{L}$.

```

1: for  $k = 0, 1, 2, \dots, K$  do
2:   for  $i = 1 \dots n$  in parallel do
3:      $x^k \leftarrow \theta_1 z^k + \theta_2 w^k + (1 - \theta_1 - \theta_2) y^k$ 
4:      $g_i^k \leftarrow \nabla f_i(x^k) - \nabla f_i(w^k)$ 
5:      $\tilde{g}_i^k \leftarrow Q_i(g_i^k)$ 
6:     Using communications broadcast  $\tilde{g}_i^k$ 
7:     Compute  $\tilde{g}^k \leftarrow \frac{1}{n} \sum_{i=1}^n \tilde{g}_i^k + \nabla f(w^k)$ 
8:      $z^{k+1} \leftarrow \frac{1}{1+\eta\sigma} (\eta\sigma x^k + z^k - \frac{\eta}{L} \tilde{g}^k)$ 
9:      $y^{k+1} \leftarrow x^k + \theta_1 (z^{k+1} - z^k)$ 
10:     $w^{k+1} \leftarrow \begin{cases} y^k, & \text{with probability } p \\ w^k, & \text{with probability } 1-p \end{cases}$ 
11:    if  $w^{k+1} = y^k$  then
12:      Using communications broadcast  $\nabla f_i(w^{k+1})$ 
13:      Compute  $\nabla f(w^{k+1}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w^{k+1})$ 
14:    end if
15:  end for
16: end for

```

DVPL-Katyusha

In DVPL-Katyusha every worker has its own set of features. Therefore, unlike the horizontal regime, in the vertical case all operations are performed on subvectors, corresponding to individual worker's components.

Algorithm 2 DVPL-Katyusha

Parameters: $\theta_1, \theta_2, \gamma \in (0, 1]$, probability $p \in (0, 1]$, RandK select j -th sample with probability $p_j = \frac{L_j}{L}$ (every worker has the same random seed for random, connected with p and p_j).

Initialization: Choose $y^0 = w^0 = z^0 \in \mathbb{R}^d$, stepsize $\eta = \frac{\theta_2 \gamma}{(1+\theta_2)\theta_1}$, set $\tilde{L} = \max\{\frac{L}{K}, L\}$, $\sigma = \frac{\sigma_*}{L}$.

```

1: for  $k = 0, 1, 2, \dots, K$  do
2:   for  $i = 1 \dots n$  in parallel do
3:      $x_i^k \leftarrow \theta_1 z_i^k + \theta_2 w_i^k + (1 - \theta_1 - \theta_2) y_i^k$ 
4:     Compute  $X_i^k = \text{RandK}\left(\left\|\langle A_{j_i}^T, x_i^k \rangle\right\|_{j=1, \dots, s}\right)$ 
5:     Compute  $W_i^k = \text{RandK}\left(\left\|\langle A_{j_i}^T, w_i^k \rangle\right\|_{j=1, \dots, s}\right)$ 
6:     Using communications broadcast  $X_i^k$  and  $W_i^k$ 
7:      $J^k = \{j_1^k, \dots, j_n^k\}$  - indices, selected by RandK
8:      $\tilde{g}_i^k \leftarrow \frac{1}{K} \sum_{j \in J^k} \frac{1}{s_j} \nabla \mathcal{L}_j\left(\sum_{i=1}^n X_{ij}^k, b_j\right) - \frac{1}{K} \sum_{j \in J^k} \frac{1}{s_j} \nabla \mathcal{L}_j\left(\sum_{i=1}^n W_{ij}^k, b_j\right) + \nabla \mathcal{L}(Aw^k, b)_i$ 
9:      $z_i^{k+1} \leftarrow \frac{1}{1+\eta\sigma} (\eta\sigma x_i^k + z_i^k - \frac{\eta}{L} \tilde{g}_i^k)$ 
10:     $y_i^{k+1} \leftarrow x_i^k + \theta_1 (z_i^{k+1} - z_i^k)$ 
11:     $w_i^{k+1} \leftarrow \begin{cases} y_i^k, & \text{with probability } p \\ w_i^k, & \text{with probability } 1-p \end{cases}$ 
12:    if  $w_i^{k+1} = y_i^k$  then
13:      for  $j = 1 \dots s$  do
14:        Compute  $\langle A_{j_i}^T, w_i^k \rangle$ 
15:        Using communications broadcast  $\langle A_{j_i}^T, w_i^k \rangle$ 
16:      end for
17:      Compute  $\nabla \mathcal{L}(Aw^k, b)_i$ 
18:    end if
19:  end for
20: end for

```

Numerical experiments

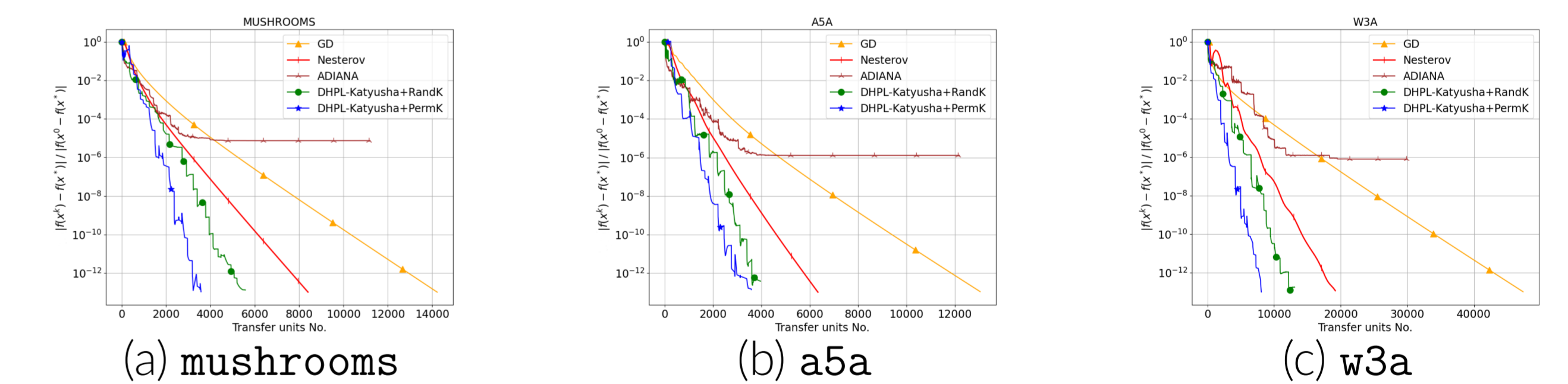


Figure 2. Comparison of different algorithms in horizontal case on LIBSVM datasets mushrooms, a5a and w3a.

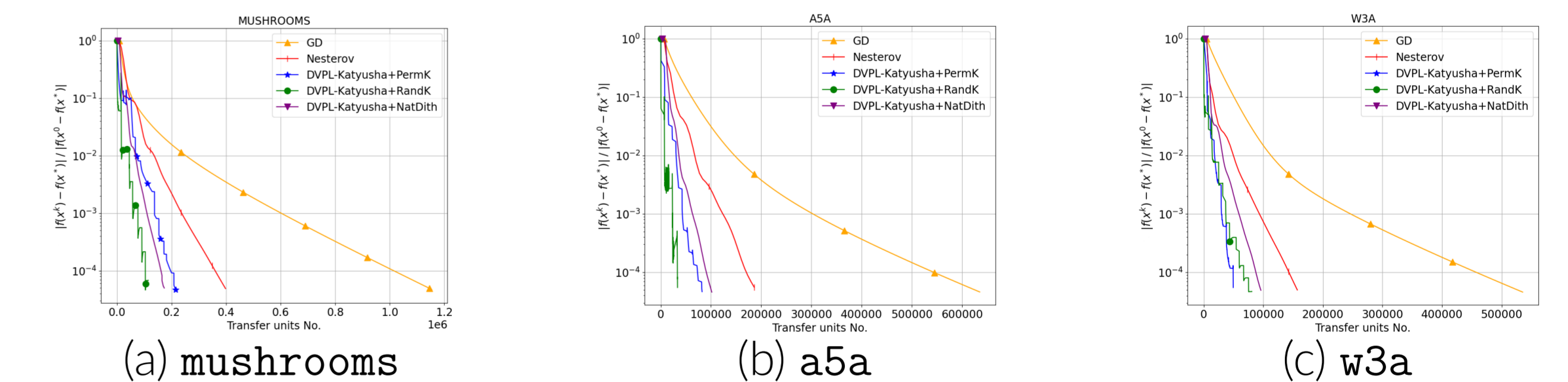


Figure 3. Comparison of different algorithms in vertical case on LIBSVM datasets mushrooms, a5a and w3a.

References

- Trishul Chilimbi, Yutaka Suzue, Johnson Apacible, and Karthik Kalyanaraman. Project adam: Building an efficient and scalable deep learning training system. In *11th USENIX symposium on operating systems design and implementation (OSDI 14)*, pages 571–582, 2014.
- Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- Dmitry Kovalev, Samuel Horvath, and Peter Richtarik. Don't jump through hoops and remove those loops: Svrq and katyusha are better without the outer loop, 2019.