

# Distributed First-Order Optimization with Tamed Communications

Dmitry Grishchenko<sup>1,2</sup> Franck IUTZELER<sup>1,2</sup> Jérôme MALICK<sup>2,3</sup>

<sup>1</sup>Université Grenoble Alpes <sup>2</sup>LJK <sup>3</sup>CNRS

## Problem

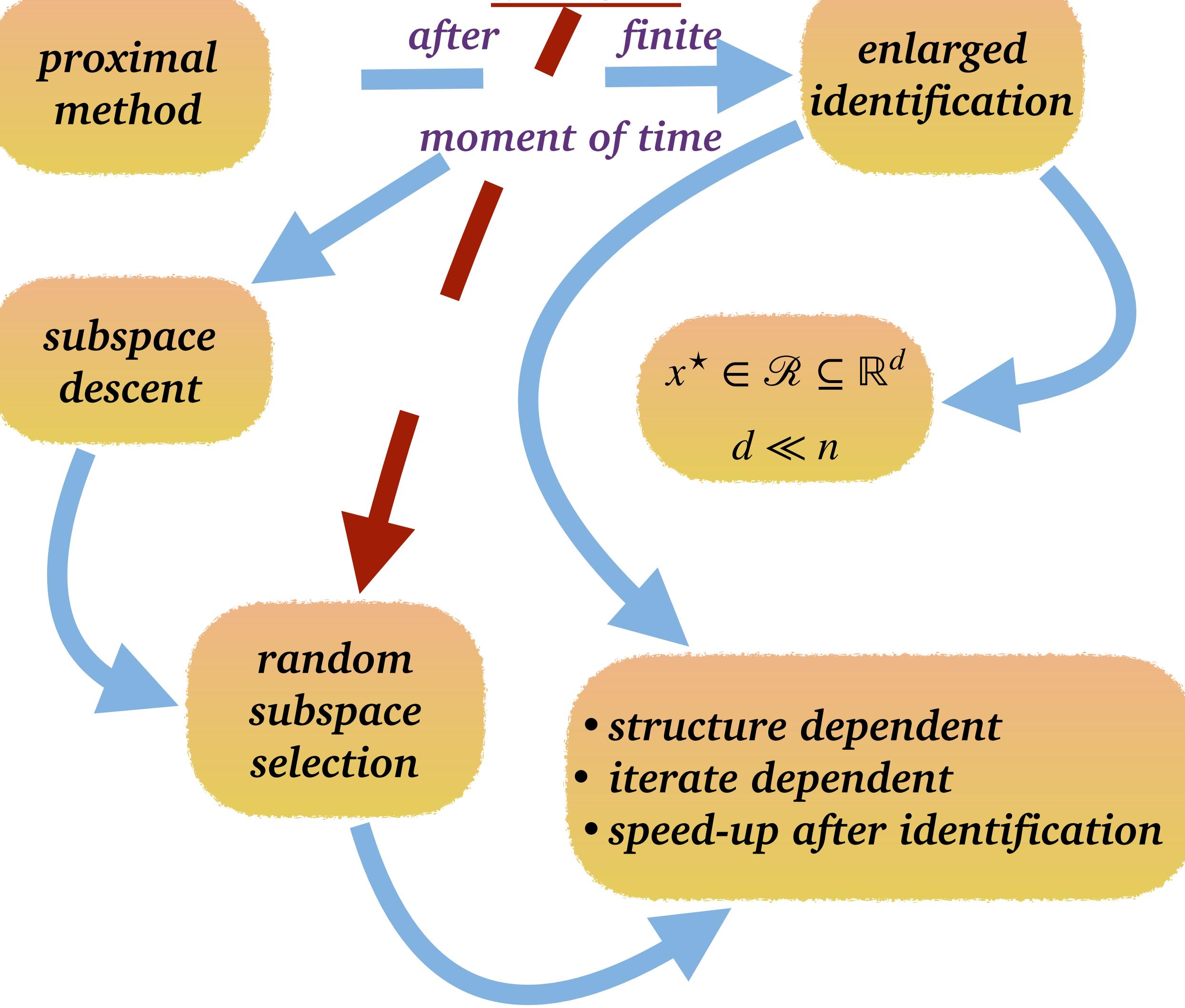
$$\min_{x \in \mathbb{R}^n} \sum_{i=1}^M \pi_i f_i(x) + g(x)$$

amount of machines  $M$  convex, proper, l.s.c.  
 $x^*$ - unique minimizer  
 BIG proportion on  $i^{th}$  machine  $L$ -smooth  $\mu$ - s. convex

## Motivation



## Intuition



Adaptive Distributed Randomized Proximal Subspace Descent - ADRPSD

- 1: [M] Generate the first admissible selection  $\mathfrak{S}^0$ , compute
- 2: [M] Initialize  $z^0, x^1 = \text{prox}_{\gamma g}(Q_0^{-1}(z^0)), \ell = 0, L = \{0\}$ . [SPARSE]
- 3: [W<sub>i</sub>] Receive  $\mathfrak{S}^0$  from master
- 4: [M, W<sub>i</sub>] Compute  $P_0 = \mathbb{E}[P_{i,\mathfrak{S}^0}]$ ,  $Q_0 = P_0^{-\frac{1}{2}}$ , and  $Q_0^{-1}$
- 5: **for**  $k = 1, \dots$  **in parallel do**
- 6: [W<sub>i</sub>] Receive  $x^k$  from master
- 7: [W<sub>i</sub>] Select independently  $P_{i,\mathfrak{S}^k}$  [SPARSE for some g]
- 8: [W<sub>i</sub>]  $y_i^k = Q_\ell(x^k - \gamma \nabla f_i(x^k))$
- 9: [W<sub>i</sub>] Send  $P_{i,\mathfrak{S}^k}(y_i^k)$  to master [SPARSE]
- 10: [M]  $z^k = \sum_{i=1}^M \pi_i(P_{i,\mathfrak{S}^k}(y_i^k) + (I - P_{i,\mathfrak{S}^k})(z^{k-1}))$
- 11: [M]  $x^{k+1} = \text{prox}_{\gamma g}(Q_\ell^{-1}(z^k))$
- 12: **if** an adaptation is decided **then**
- 13: [M]  $L \leftarrow L \cup \{k+1\}, \ell \leftarrow \ell + 1$
- 14: [M] Generate a new admissible selection  $\mathfrak{S}^\ell$  [SPARSE]
- 15: [W<sub>i</sub>] Receive  $\mathfrak{S}^\ell$  from master
- 16: [M, W<sub>i</sub>] Compute  $P_\ell = \mathbb{E}[P_{i,\mathfrak{S}^\ell}]$ ,  $Q_\ell = P_\ell^{-\frac{1}{2}}$ , and  $Q_\ell^{-1}$
- 17: [M] Rescale  $z^k \leftarrow Q_\ell Q_{\ell-1}^{-1} z^k$
- 18: **end if**
- 19: **end for**

## Subspace families and projections

The family of linear subspaces  $\mathcal{C} = \{\mathcal{C}_i\}_i$  is called covering if  $\sum_i \mathcal{C}_i = \mathbb{R}^n$

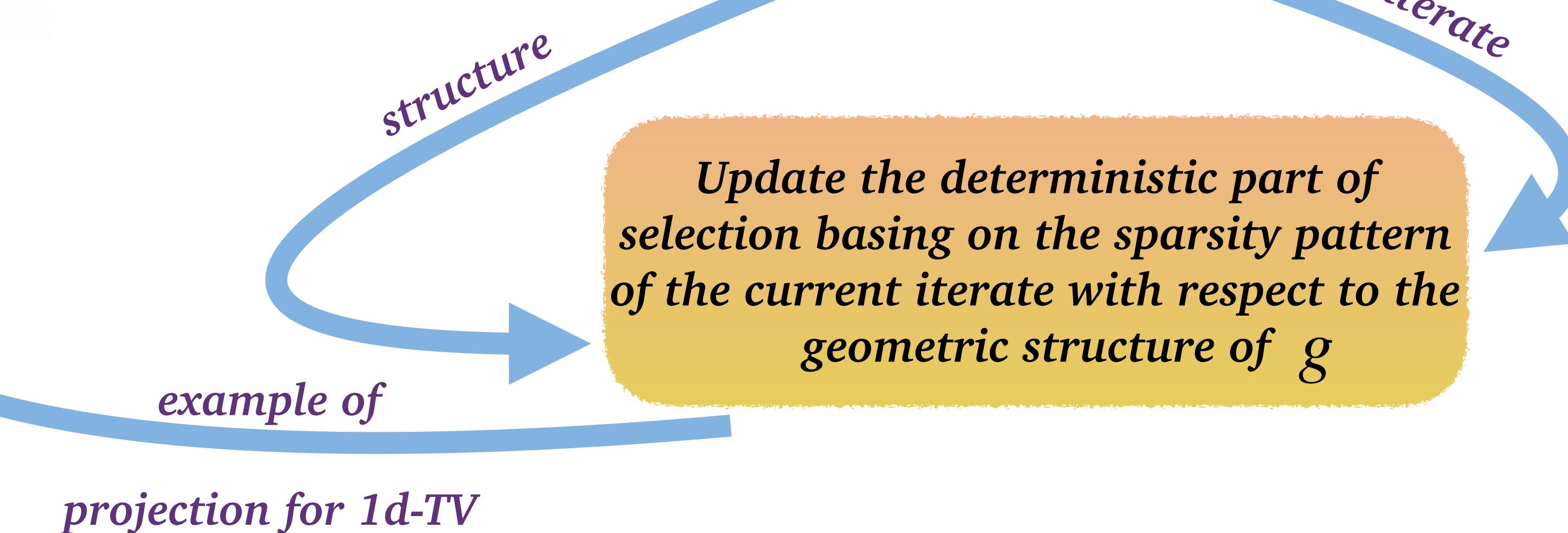
The random generation of a subspace by selecting randomly some subspaces in some covering family is called admissible selection if it samples whole space.

$$P_{\mathfrak{S}} = \begin{Bmatrix} \overbrace{\frac{1}{n_1} \dots \frac{1}{n_1}}^{n_1} & 0 & \dots & \overbrace{\dots \dots \dots}^{n-n_s} & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \frac{1}{n_1} \dots \frac{1}{n_1} & 0 & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & 0 & \frac{1}{n-n_s} \dots \frac{1}{n-n_s} \\ 0 & \dots & \dots & 0 & \frac{1}{n-n_s} \dots \frac{1}{n-n_s} \end{Bmatrix}_{n-n_s}$$

If  $\mathfrak{S}$  is admissible selection, then (average projection) is invertible.  $P := \mathbb{E}[P_{\mathfrak{S}}]$

$P_{\mathfrak{F}}$  is an orthogonal projection onto linear subspace  $\mathfrak{F}$

## Adaptive selection



**ADRPSD convergence** For any  $\gamma \in (0, 2/(\mu + L)]$ , let the user choose its adaptation strategy so that:

- the *adaptation cost* is upper bounded by a deterministic sequence:  $\|Q_\ell Q_{\ell-1}^{-1}\|_2^2 \leq \mathbf{a}_\ell$ ;
- the *inter-adaptation time* is lower bounded by a deterministic sequence:  $k_\ell - k_{\ell-1} \geq \mathbf{c}_\ell$ ;
- the *selection uniformity* is lower bounded by a deterministic sequence:  $\lambda_{\min}(P_\ell) \geq \lambda_\ell$ ;

then, from the previous instantaneous rate  $1 - \alpha_{\ell-1} := 1 - 2\gamma\mu L\lambda_{\ell-1}/(\mu + L)$ , the corrected rate for cycle  $\ell$  writes

$$(1 - \beta_\ell) := (1 - \alpha_{\ell-1})\mathbf{a}_\ell^{1/c_\ell}.$$

Then, we have for any  $k \in [k_\ell, k_{\ell+1}]$

$$\mathbb{E}[\|x^{k+1} - x^*\|_2^2] \leq (1 - \alpha_\ell)^{k-k_\ell} \prod_{m=1}^{\ell} (1 - \beta_m)^{c_m} \|z^0 - Q_0(x^* - \gamma \nabla f(x^*))\|_2^2.$$

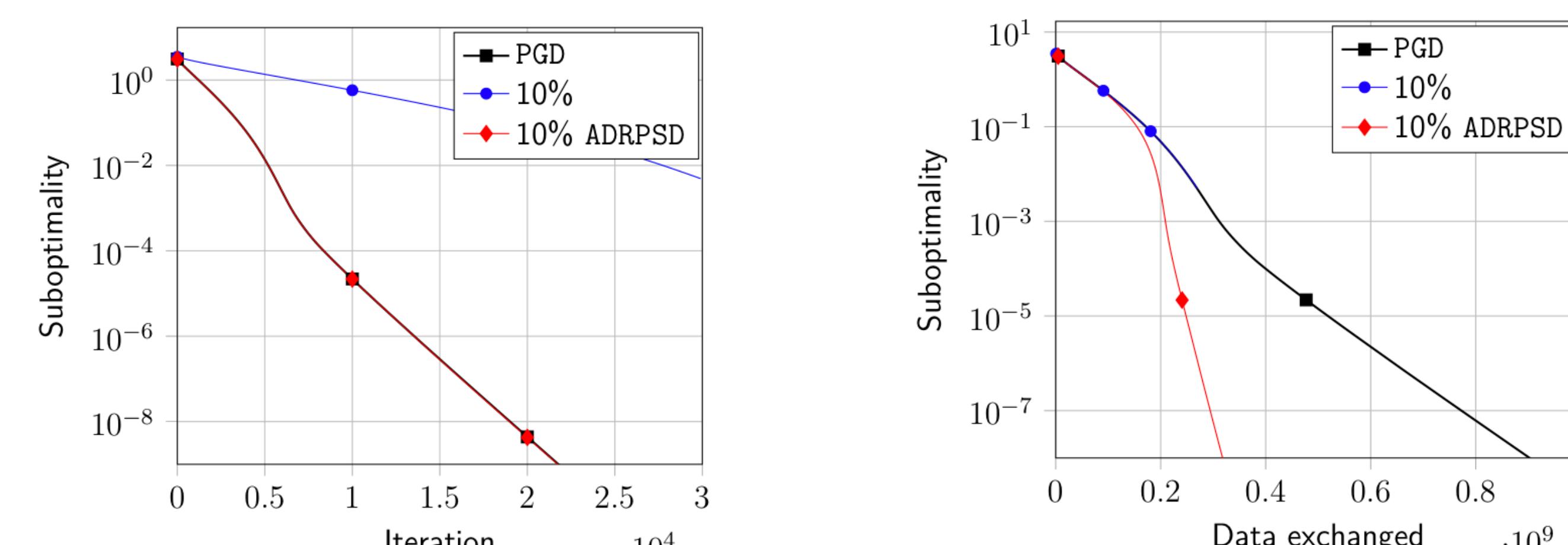


Figure 1:  $\ell_1$  regularized logistic regression on rcv\_1 dataset

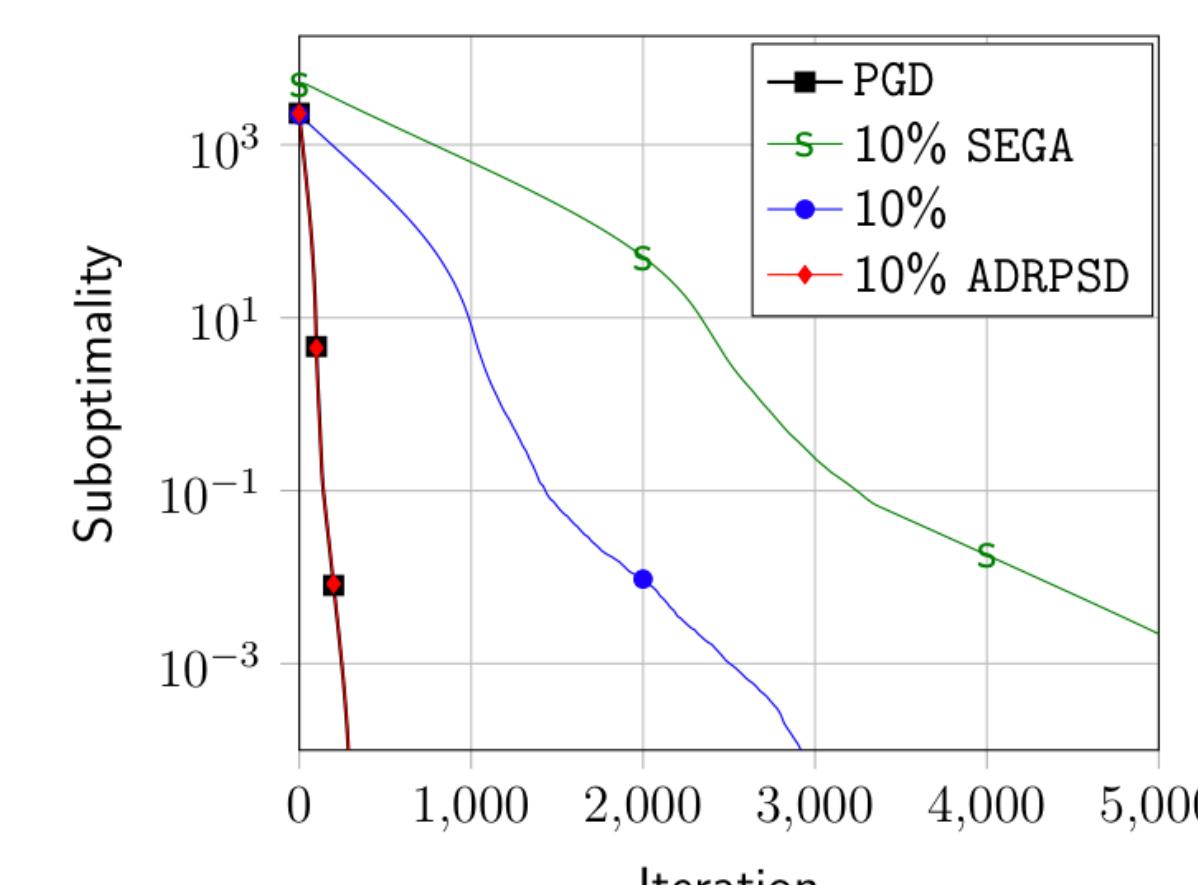


Figure 2:  $\ell_{1,2}$  regularized logistic regression on rcv\_1 dataset

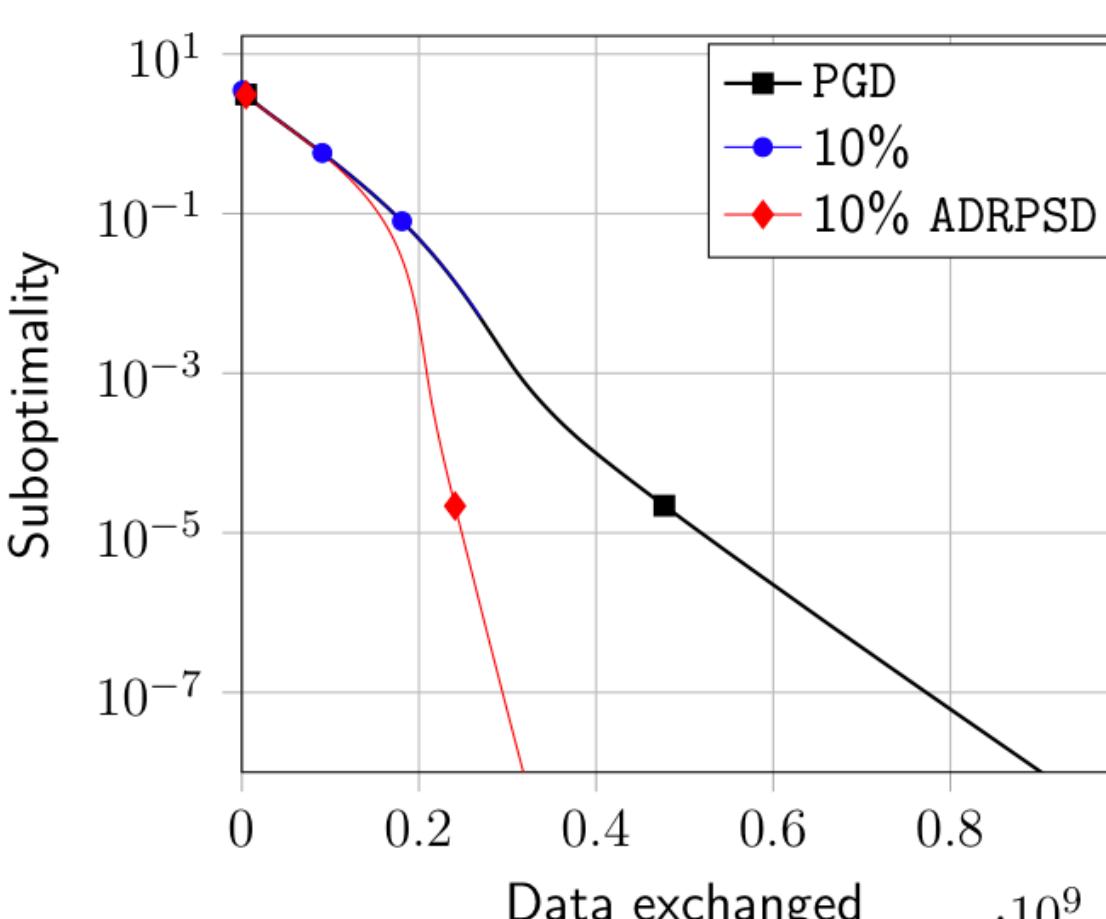
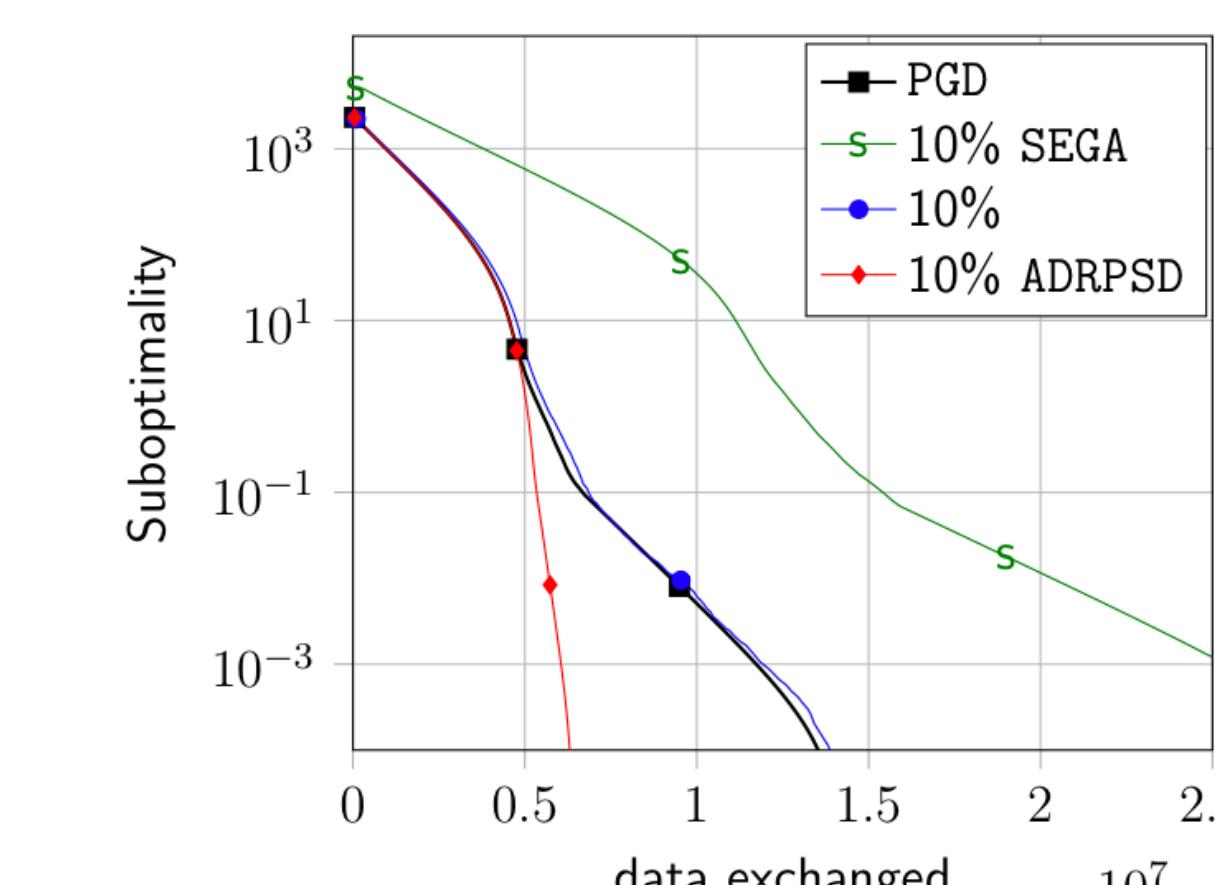


Figure 3: TV regularized logistic regression on a1a dataset



- [1] Bach, F., Jenatton, R., Mairal, J., Obozinski, G.: Optimization with sparsity-inducing penalties. Foundations and Trends® in Machine Learning (2012)  
 [2] Fadili, J., Malick, J., Peyré, G.: Sensitivity analysis for mirror-stratifiable convex functions. SIAM Journal on Optimization (2018)  
 [3] Grishchenko, D., et al.: Asynchronous distributed learning with sparse communications and identification. arXiv preprint arXiv:1812.03871 (2018)  
 [4] Hanzely, F., Mishchenko, K., Richtárik, P.: SegA: Variance reduction via gradient sketching. In: Advances in NeurIPS (2018)