

# Sparse Asynchronous Distributed Learning

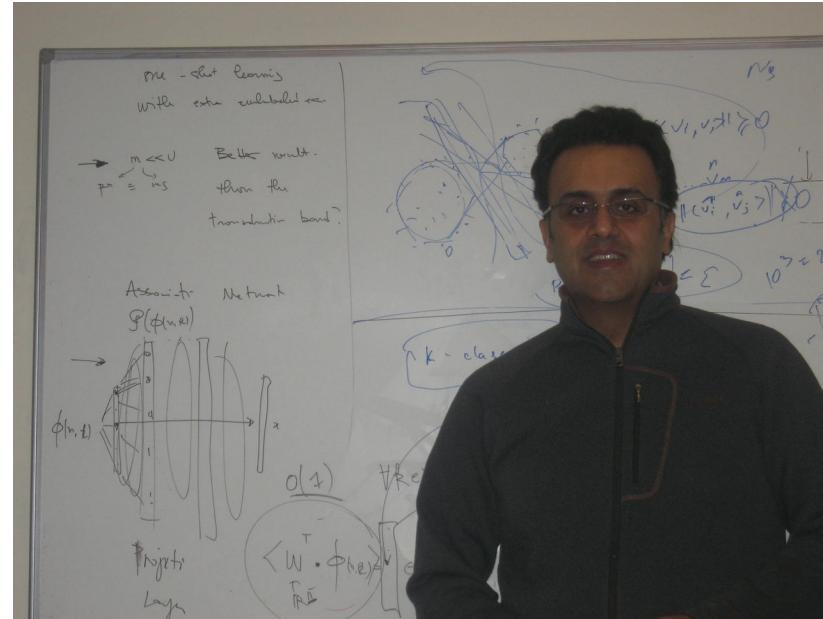
*Dmitry Grishchenko*

18 November 2020

# Collaborators



**Frank Lutzeler**



**Massih-Reza Amini**

# Problem

$$\min_{x \in \mathbb{R}^n} \underbrace{\frac{1}{m} \sum_{i=1}^m \ell(b_i, h(a_i, x)) + \lambda_1 \|x\|_1}_{f(x)}$$

# Problem

$$\min_{x \in \mathbb{R}^n} \sum_{i=1}^M \alpha_i \underbrace{\left[ \frac{1}{|\mathcal{D}_i|} \sum_{j \in \mathcal{D}_i} \ell(b_j, h(a_j, x)) \right]}_{f_i} + \lambda_1 \|x\|_1,$$

where the full dataset  $\mathcal{D}$  is split onto  $M$  nonintersecting subsets  $\mathcal{D}_i$  and  $\alpha_i$  is the proportion of examples  $\frac{|\mathcal{D}_i|}{m}$ .

# Problem

$$\min_{x \in \mathbb{R}^n} \sum_{i=1}^M \alpha_i \underbrace{\left[ \frac{1}{|\mathcal{D}_i|} \sum_{j \in \mathcal{D}_i} \ell(b_j, h(a_j, x)) \right]}_{f_i} + \lambda_1 \|x\|_1,$$

$L$ -smooth  
 $\mu$ -strongly convex

where the full dataset  $\mathcal{D}$  is split onto  $M$  nonintersecting subsets  $\mathcal{D}_i$  and  $\alpha_i$  is the proportion of examples  $\frac{|\mathcal{D}_i|}{m}$ .

# Problem

$$\min_{x \in \mathbb{R}^n} \sum_{i=1}^M \alpha_i \underbrace{\left[ \frac{1}{|\mathcal{D}_i|} \sum_{j \in \mathcal{D}_i} \ell(b_j, h(a_j, x)) \right]}_{f_i} + \lambda_1 \|x\|_1,$$

*L-smooth*  
 *$\mu$ -strongly convex*

where the full dataset  $\mathcal{D}$  is split onto  $M$  nonintersecting subsets  $\mathcal{D}_i$  and  $\alpha_i$  is the proportion of examples  $\frac{|\mathcal{D}_i|}{m}$ .

These subsets  $\mathcal{D}_i$  can be split over machines.



**Worker**

**Communication**



**Master**

# Problem

$$\min_{x \in \mathbb{R}^n} \sum_{i=1}^M \alpha_i \underbrace{\left[ \frac{1}{|\mathcal{D}_i|} \sum_{j \in \mathcal{D}_i} \ell(b_j, h(a_j, x)) \right]}_{f_i} + \lambda_1 \|x\|_1,$$

*L-smooth*  
 *$\mu$ -strongly convex*

where the full dataset  $\mathcal{D}$  is split onto  $M$  nonintersecting subsets  $\mathcal{D}_i$  and  $\alpha_i$  is the proportion of examples  $\frac{|\mathcal{D}_i|}{m}$ .

These subsets  $\mathcal{D}_i$  can be split over machines.



**Worker**

**Communication**

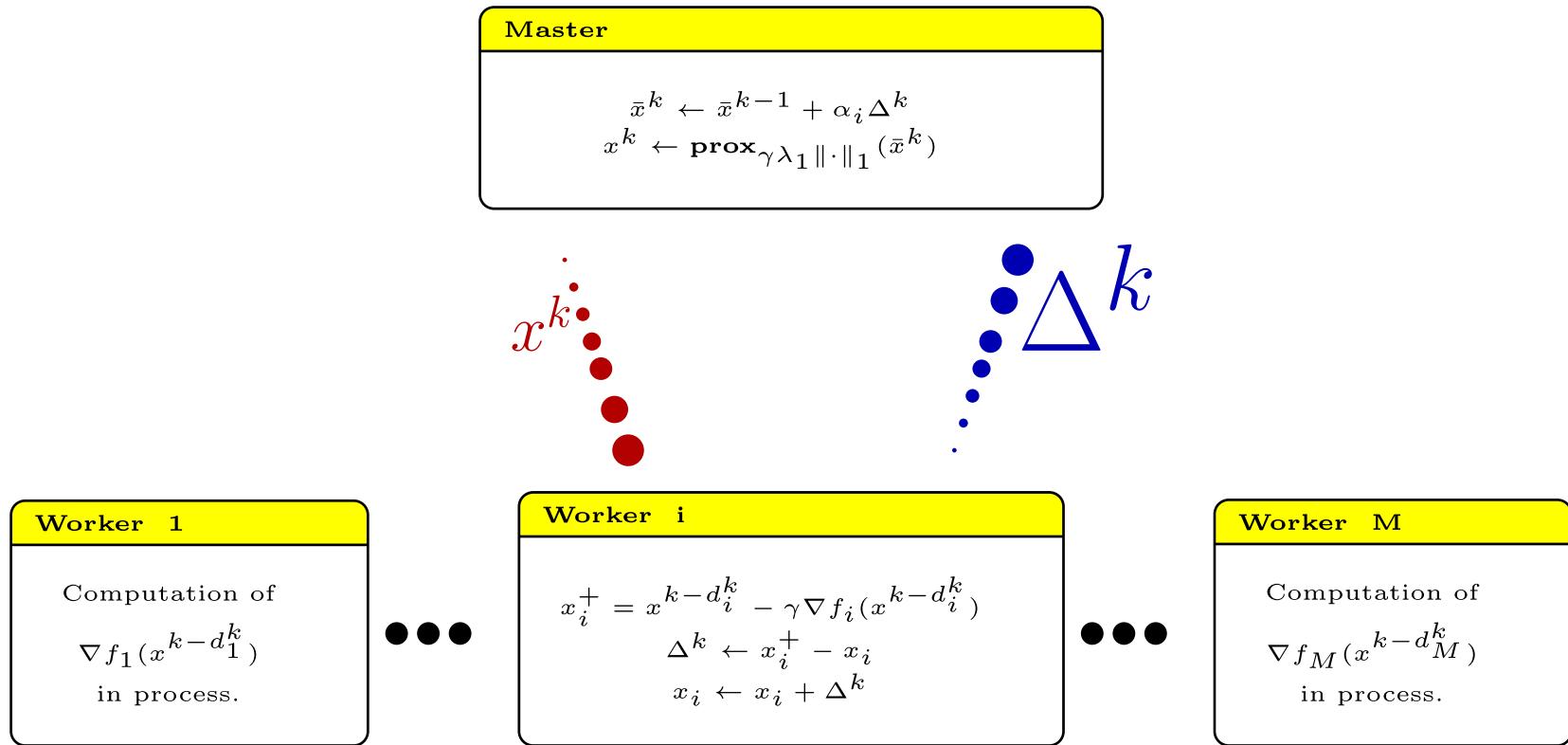


**Bottleneck**



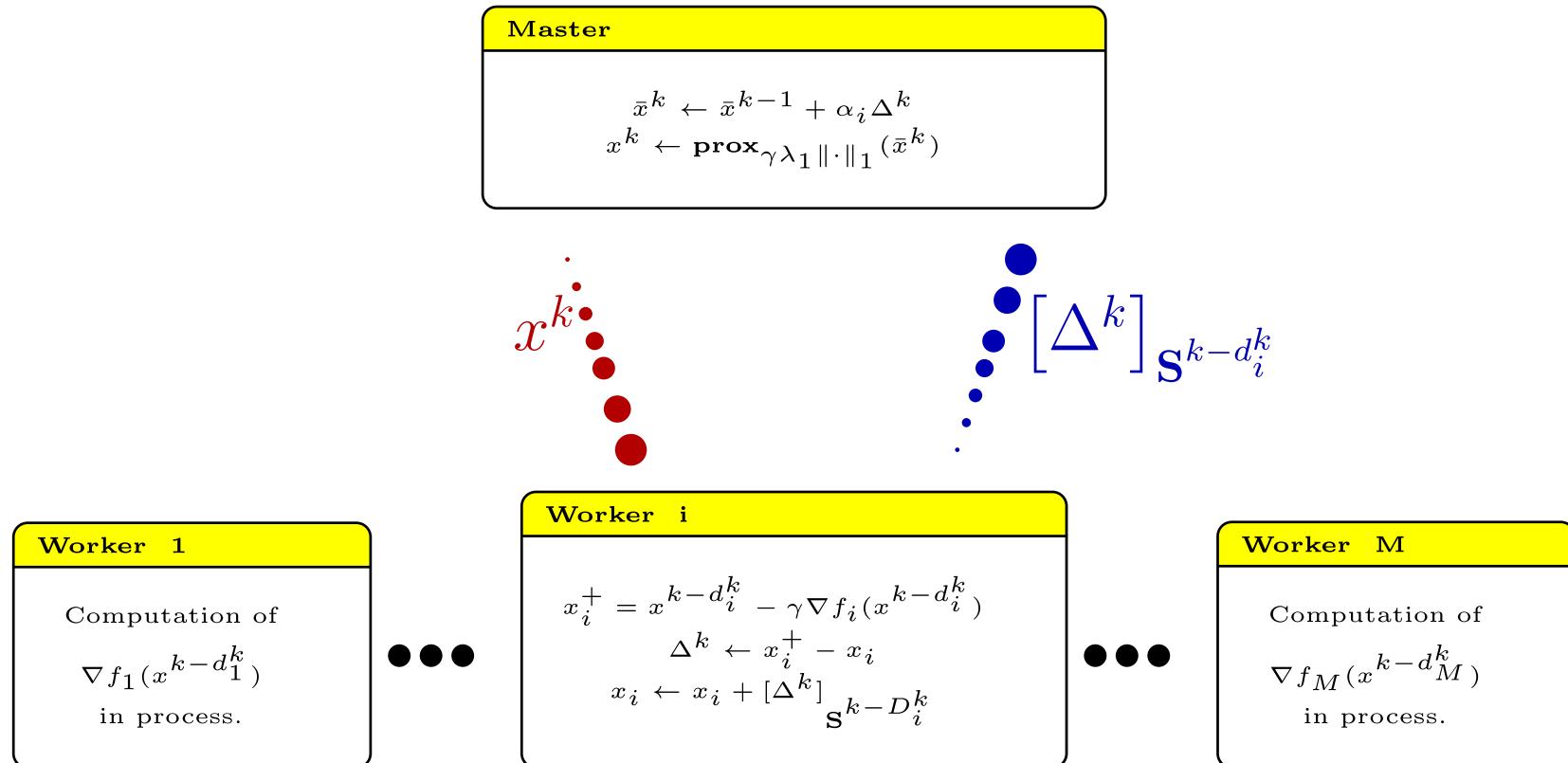
**Master**

# Algorithm: DAve-PG

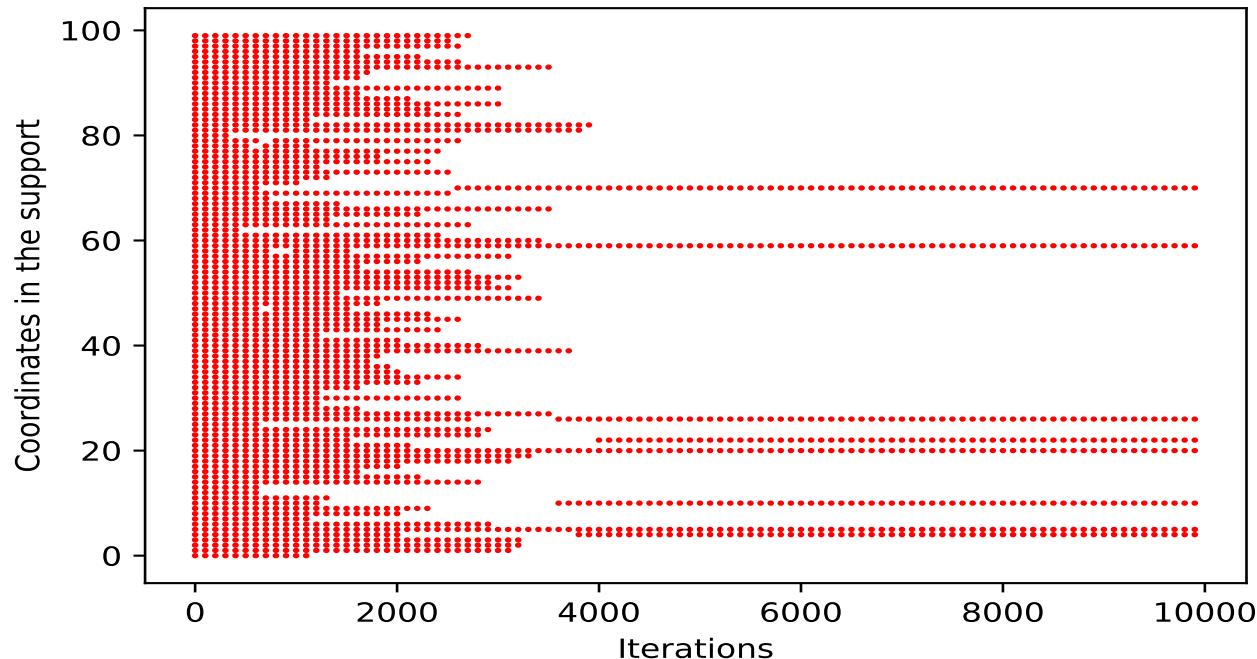


Konstantin Mishchenko, Franck Iutzeler, Jérôme Malick, and Massih-Reza Amini. *A Delay-tolerant Proximal-Gradient Algorithm for Distributed Learning*, International Conference on Machine Learning, 3584-3592

# Algorithm: SPY



# Identification (Example)



Synthetic LASSO problem  $\min \frac{1}{2} \|Ax - b\|_2^2 + \lambda_1 \|x\|_1$  for random generated matrix  $A \in \mathbb{R}^{100 \times 100}$  and vector  $b \in \mathbb{R}^{100}$  and hyperparameter  $\lambda_1$  chosen to reach 8% of density (amount of non-zero coordinates) of the final solution.

# Adaptive Selection

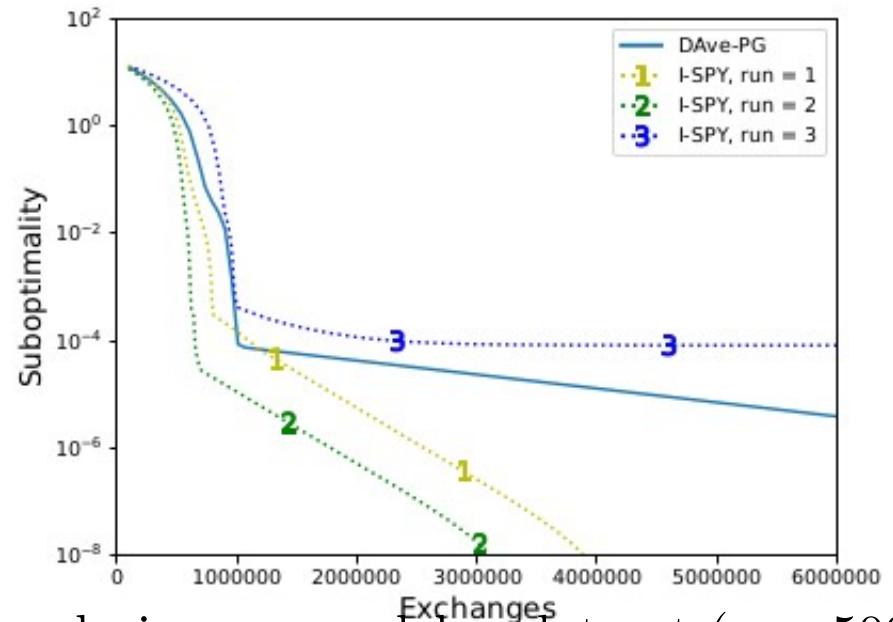
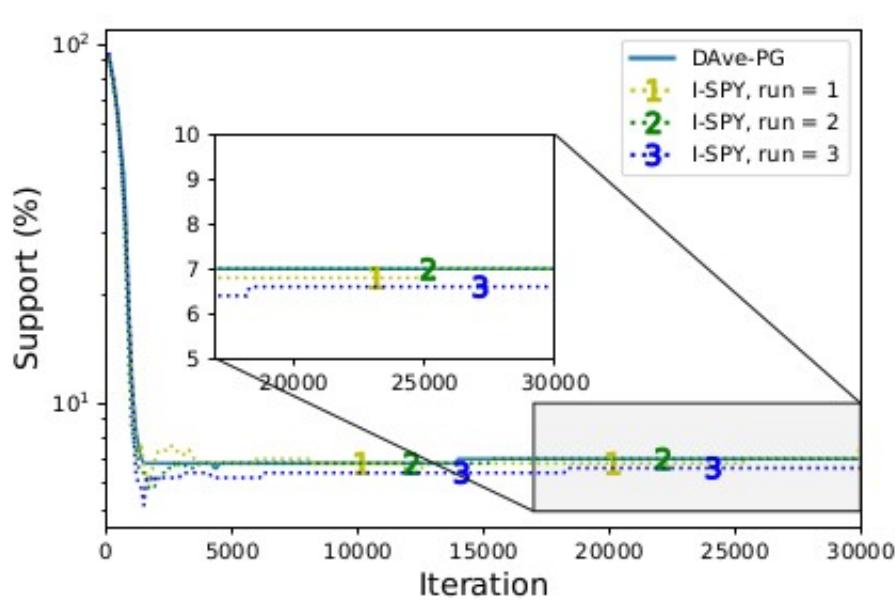
$p$  is  $\pi$ -priority random vector w.r.t. the current iterate point  $x^k$

$$\mathbb{P} [j \in \mathbf{S}_\pi^k] = \begin{cases} 1 & \text{if } j \in \text{supp}(x^k), \\ \pi & \text{otherwise.} \end{cases}$$

**This selection is not i.i.d.!**

**If support is fixed the selection is i.i.d.!**

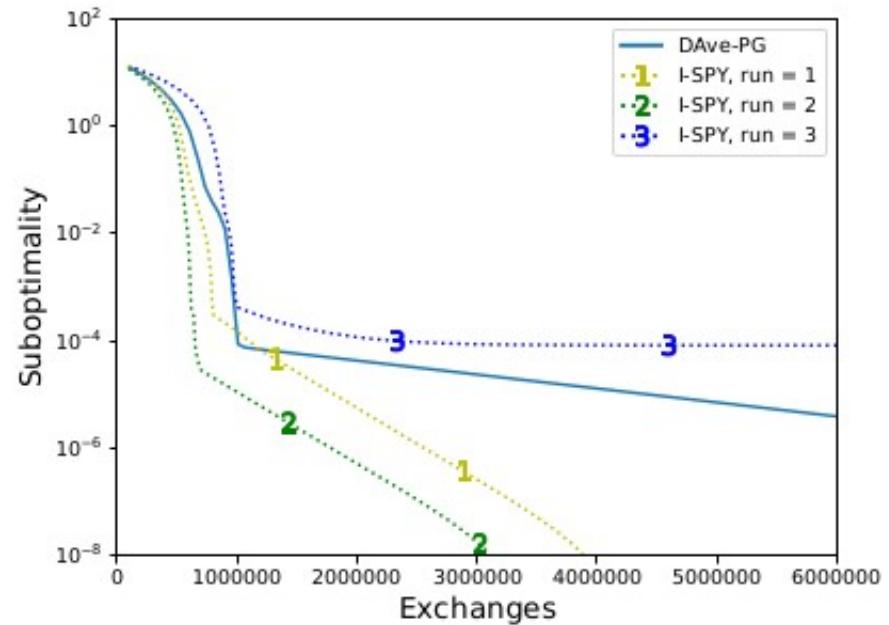
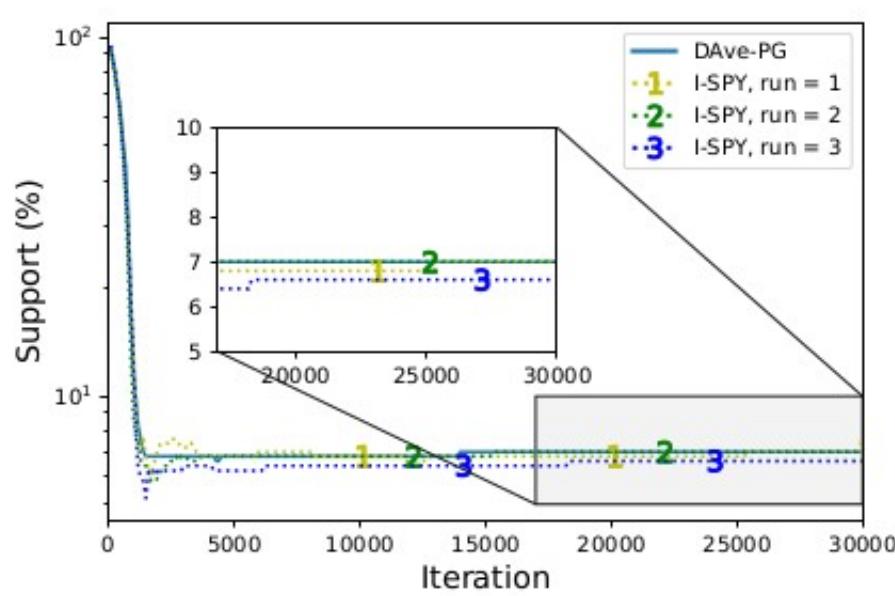
# Experiments: Adaptive Selection



Logistic regression with elastic net regularizer on madelon dataset ( $n = 500$ ,  $m = 2000$ ) and  $M = 10$  machines.

$$\min_{x \in \mathbb{R}^n} \frac{1}{m} \sum_{j=1}^m \log(1 + \exp(-y_j z_j^\top x)) + \lambda_1 \|x\|_1 + \frac{\lambda_2}{2} \|x\|_2^2$$

# Experiments: Adaptive Selection



It is better if it converges, but it can diverge!

# Mask Selection

# Mask Selection

$$[\Delta^k]_{\mathbf{S}^k}$$

# Mask Selection

$$[\Delta^k]_{\mathbf{S}^k}$$

Random sparsification with  $p = (p_1, \dots, p_n) \in (0, 1]^n$ .

$$\mathbb{P}[j \in \mathbf{S}_p^k] = p_j > 0 \quad \text{for all } j \in \{1, \dots, n\}.$$

# Mask Selection

$$[\Delta^k]_{\mathbf{S}^k}$$

Random sparsification with  $p = (p_1, \dots, p_n) \in (0, 1]^n$ .

$$\mathbb{P}[j \in \mathbf{S}_p^k] = p_j > 0 \quad \text{for all } j \in \{1, \dots, n\}.$$

- $p$  is an arbitrary probability vector.

# Mask Selection

$$[\Delta^k]_{\mathbf{S}^k}$$

Random sparsification with  $p = (p_1, \dots, p_n) \in (0, 1]^n$ .

$$\mathbb{P}[j \in \mathbf{S}_p^k] = p_j > 0 \quad \text{for all } j \in \{1, \dots, n\}.$$

- $p$  is an arbitrary probability vector.
- $p$  is a  $\pi$ -uniform probability vector.

# Mask Selection

$$[\Delta^k]_{\mathbf{S}^k}$$

Random sparsification with  $p = (p_1, \dots, p_n) \in (0, 1]^n$ .

$$\mathbb{P}[j \in \mathbf{S}_p^k] = p_j > 0 \quad \text{for all } j \in \{1, \dots, n\}.$$

- $p$  is an arbitrary probability vector.
- $p$  is a  $\pi$ -uniform probability vector.
- $p$  is a  $\pi$ -priority random vector w.r.t. some point  $x$

$$\mathbb{P}[j \in \mathbf{S}_\pi^k] = \begin{cases} 1 & \text{if } j \in \text{supp}(x), \\ \pi & \text{otherwise.} \end{cases}$$

# General Theoretical Result

# General Theoretical Result

## Assumption (on randomness)

The sparsity mask selectors ( $\mathbf{S}_p^k$ ) are independent and identically distributed random variables. We select a coordinate in the mask as follows:

$$\mathbb{P}[j \in \mathbf{S}_p^k] = p_j > 0 \quad \text{for all } j \in \{1, \dots, n\},$$

with  $p = (p_1, \dots, p_n) \in (0, 1]^n$ .

# General Theoretical Result

## Assumption (on randomness)

The sparsity mask selectors ( $\mathbf{S}_p^k$ ) are independent and identically distributed random variables. We select a coordinate in the mask as follows:

$$\mathbb{P}[j \in \mathbf{S}_p^k] = p_j > 0 \quad \text{for all } j \in \{1, \dots, n\},$$

with  $p = (p_1, \dots, p_n) \in (0, 1]^n$ .

## Theorem (Limits of sparsification)

Take  $\gamma = \frac{2}{\mu+L}$ , then SPY verifies for all  $k \in [k_m, k_{m+1})$

$$\mathbb{E} \|x^k - x^\star\|^2 \leq \left( p_{\max} \left( \frac{1 - \kappa_P}{1 + \kappa_P} \right)^2 + 1 - p_{\min} \right)^m \max_i \|x_i^0 - x_i^\star\|^2.$$

with the shifted local solutions  $x_i^\star = x^\star - \gamma_i \nabla f_i(x^\star)$ .

# General Theoretical Result

## Assumption (on randomness)

The sparsity mask selectors ( $\mathbf{S}_p^k$ ) are independent and identically distributed random variables. We select a coordinate in the mask as follows:

$$\mathbb{P}[j \in \mathbf{S}_p^k] = p_j > 0 \quad \text{for all } j \in \{1, \dots, n\},$$

with  $p = (p_1, \dots, p_n) \in (0, 1]^n$ .

## Theorem (Limits of sparsification)

Take  $\gamma = \frac{2}{\mu+L}$ , then SPY verifies for all  $k \in [k_m, k_{m+1})$

$$\mathbb{E} \|x^k - x^*\|^2 \leq \left( p_{\max} \left( \frac{1 - \kappa_P}{1 + \kappa_P} \right)^2 + 1 - p_{\min} \right)^m \max_i \|x_i^0 - x_i^*\|^2.$$

$\in (0, 1)$

with the shifted local solutions  $x_i^* = x^* - \gamma_i \nabla f_i(x^*)$ .

# General Theoretical Result

## Assumption (on randomness)

The sparsity mask selectors ( $\mathbf{S}_p^k$ ) are independent and identically distributed random variables. We select a coordinate in the mask as follows:

$$\mathbb{P}[j \in \mathbf{S}_p^k] = p_j > 0 \quad \text{for all } j \in \{1, \dots, n\},$$

with  $p = (p_1, \dots, p_n) \in (0, 1]^n$ .

## Theorem (Limits of sparsification)

Take  $\gamma = \frac{2}{\mu+L}$ , then SPY with  $\pi$ -uniform sampling verifies for all  $k \in [k_m, k_{m+1})$

$$\mathbb{E} \|x^k - x^\star\|^2 \leq \left(1 - \frac{4\mu L}{(\mu + L)^2}\right)^m \max_i \|x_i^0 - x_i^\star\|^2.$$

with the shifted local solutions  $x_i^\star = x^\star - \gamma_i \nabla f_i(x^\star)$ .

# General Theoretical Result

## Assumption (on randomness)

The sparsity mask selectors ( $\mathbf{S}_p^k$ ) are independent and identically distributed random variables. We select a coordinate in the mask as follows:

$$\mathbb{P}[j \in \mathbf{S}_p^k] = p_j > 0 \quad \text{for all } j \in \{1, \dots, n\},$$

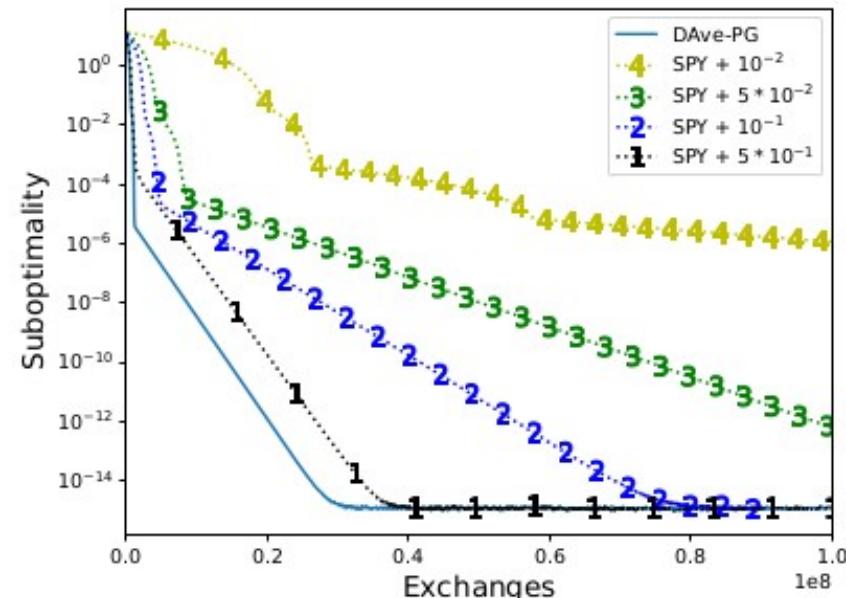
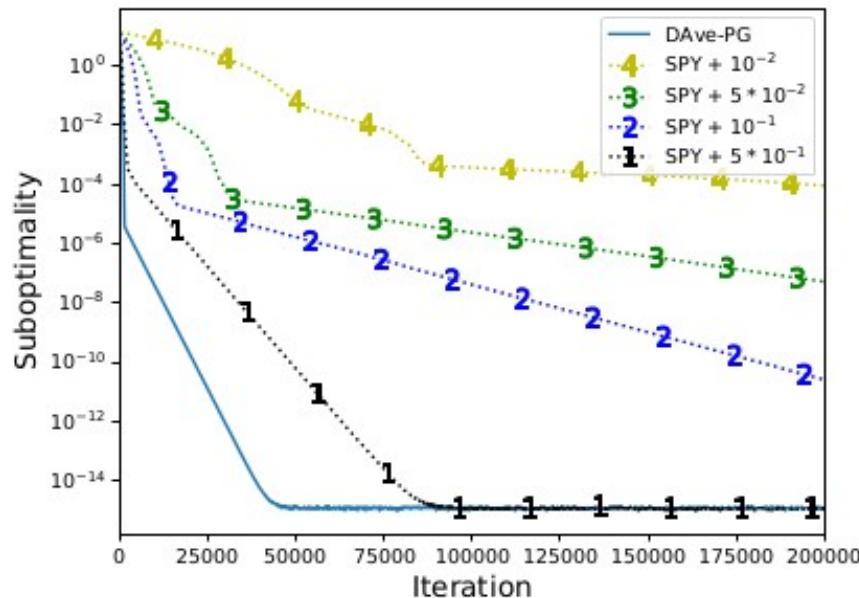
with  $p = (p_1, \dots, p_n) \in (0, 1]^n$ .

## Limits of sparsification

SPY reaches linear convergence of the mean squared error in terms of epochs if

$$\frac{p_{\min}}{p_{\max}} > (1 - \gamma\mu)^2 \stackrel{\gamma=\frac{2}{\mu+L}}{\geq} \left( \frac{1 - \kappa_P}{1 + \kappa_P} \right)^2.$$

# Uniform Sampling



Logistic regression with elastic net regularizer on madelon dataset ( $n = 500$ ,  $m = 2000$ ) and  $M = 10$  machines.

$$\min_{x \in \mathbb{R}^n} \quad \frac{1}{m} \sum_{j=1}^m \log(1 + \exp(-y_j z_j^\top x)) + \lambda_1 \|x\|_1 + \frac{\lambda_2}{2} \|x\|_2^2$$

# Identification (Example)

$$x^* = \operatorname{argmin}_{x \in \mathbb{R}^n} \sum_{i=1}^M \alpha_i f_i(x) + \lambda_1 \|x\|_1$$

## Theorem (Enlarged identification)

Let  $(u^k)$  be an  $\mathbb{R}^n$ -valued sequence converging almost surely to  $u^*$  and define sequence  $(x^k)$  as  $x^k = \mathbf{prox}_{\gamma r}(u^k)$  and  $x^* = \mathbf{prox}_{\gamma r}(u^*)$ . Then  $(x^k)$  identifies some subspaces with probability one; more precisely for any  $\varepsilon > 0$ , with probability one, after some finite time,

$$\text{supp}(x^*) \subseteq \text{supp}(x^k) \subseteq \max_{u \in \mathcal{B}(u^*, \varepsilon)} \{\text{supp}(\mathbf{prox}_{\gamma r}(u))\}.$$

# Non-degeneracy

Another way to define the non-degeneracy for the problem

$$\min_{x \in \mathbb{R}^n} f(x) + r(x)$$

is the following:

$$\nabla f(x^\star) \in \text{ri } \partial r(x^\star).$$

In case of  $\ell_1$  regularizer  $r(x) = \lambda_1 \|x\|_1$  this can be written explicitly as

$$|\nabla f(x^\star)_{[j]}| < \lambda_1 \quad \text{for all } j \in \text{supp}(x^\star).$$

# Better rate

## Assumption (Convergence)

Let us assume that for any  $\varepsilon > 0$  there exists iterate number  $K$  such that for any  $k > K$ , the average point  $\|\bar{x}^k - \bar{x}^\star\|_2^2 < \varepsilon$  is  $\varepsilon$ -close to the final solution.

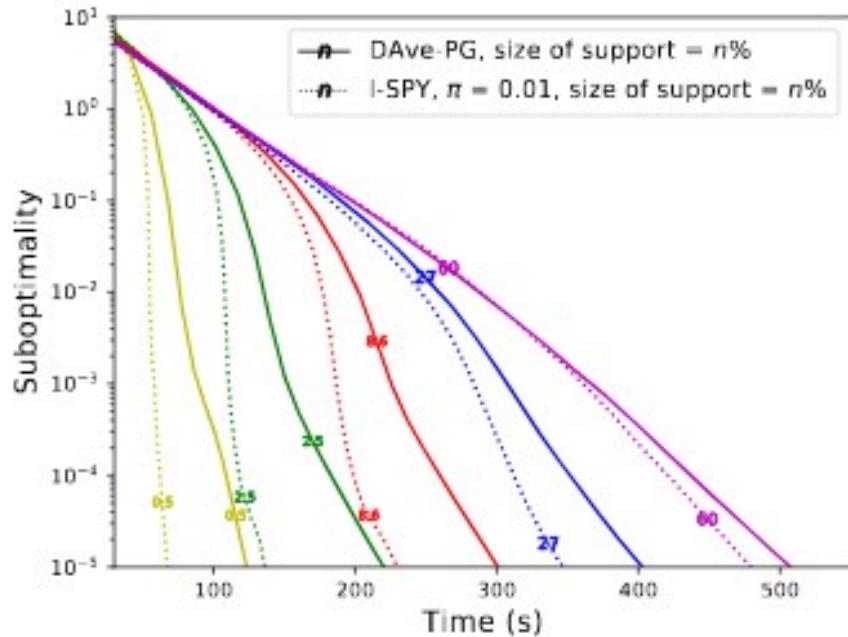
## Theorem (Better rate)

Suppose that Assumption holds. For any  $\gamma \in (0, 2/(\mu + L)]$  and for any  $k \in [k_s, k_{s+1})$  we have:

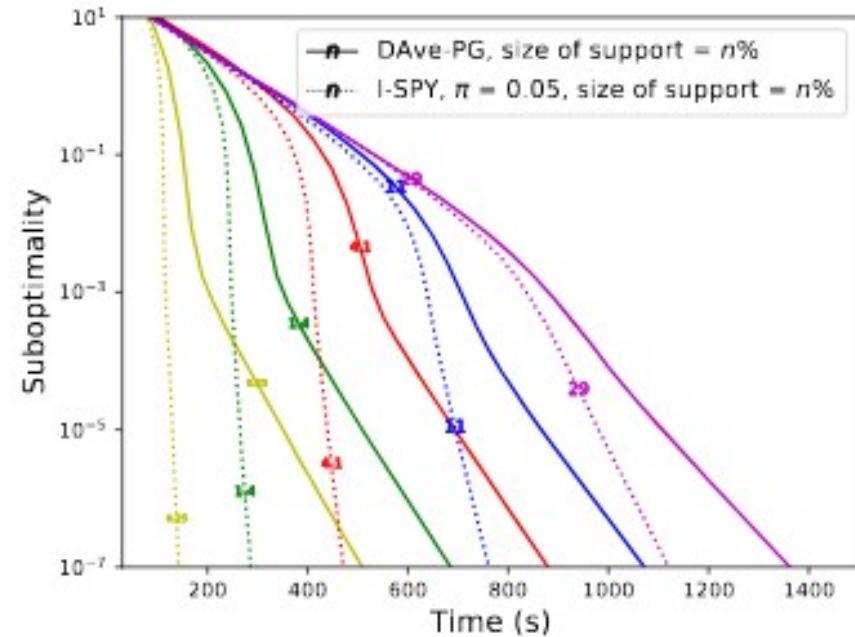
$$\|x^k - x^\star\|^2 = \mathcal{O}_p \left( \left( 1 - \frac{2\gamma\mu L}{\mu + L} \right)^s \right),$$

where  $\mathcal{O}_p$  denotes big  $O$  in probability.

# Time performance

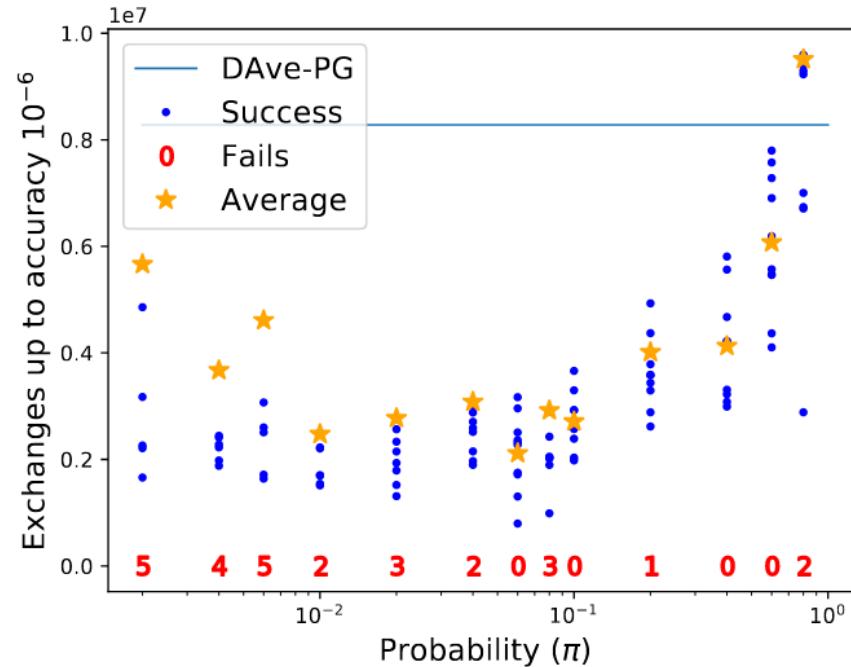


(a) real-sim



(b) rcv1\_train

# Stability



Logistic regression with elastic net regularizer on madelon dataset ( $n = 500$ ,  $m = 2000$ ) and  $M = 10$  machines.

$$\min_{x \in \mathbb{R}^n} \frac{1}{m} \sum_{j=1}^m \log(1 + \exp(-y_j z_j^\top x)) + \lambda_1 \|x\|_1 + \frac{\lambda_2}{2} \|x\|_2^2$$



***Thank You For***

***Your Attention!***