

Distributed Optimization with Flexible Communications

Franck Iutzeler LJK, Univ. Grenoble Alpes

PGMO Days 2020



Composite minimization

$$\begin{array}{ll} \text{Find } x^* \in \arg \min_{x \in \mathbb{R}^n} & \mathcal{R}(x; \{a_i, b_i\}_{i=1}^m) + \lambda r(x) \\ \text{Find } x^* \in \arg \min_{x \in \mathbb{R}^n} & f(x) + g(x) \\ & \text{smooth} \qquad \qquad \text{non-smooth} \end{array}$$

- > f : differentiable surrogate of the empirical risk \Rightarrow **Gradient**
non-linear smooth function that depends on all the data
- > g : non-smooth but chosen regularization \Rightarrow **Proximity operator**
non-differentiability on some manifolds implies structure on the solutions

closed form/easy for many regularizations:

$$\mathbf{prox}_{\gamma g}(u) = \arg \min_{y \in \mathbb{R}^n} \left\{ g(y) + \frac{1}{2\gamma} \|y - u\|_2^2 \right\}$$

- $g(x) = \|x\|_1$
- $g(x) = TV(x)$
- $g(x) = \text{indicator}_C(x)$

Natural optimization method: **proximal gradient**

$$\begin{cases} u_{k+1} = x_k - \gamma \nabla f(x_k) \\ x_{k+1} = \mathbf{prox}_{\gamma g}(u_{k+1}) \end{cases}$$

and its stochastic variants: proximal sgd, etc.

Example: LASSO

$$\begin{aligned} \text{Find } x^* \in \arg \min_{x \in \mathbb{R}^n} \mathcal{R}(x; \{a_i, b_i\}_{i=1}^m) &+ \lambda r(x) \\ \text{Find } x^* \in \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2 &+ \lambda \|x\|_1 \\ \text{smooth} & \qquad \qquad \qquad \text{non-smooth} \end{aligned}$$

Coordinates **Structure** \leftrightarrow **Optimality conditions**

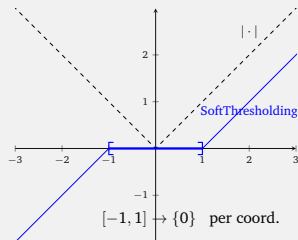
$$\forall i \quad x_i^* = 0 \quad \Leftrightarrow \quad A_i^\top (Ax^* - b) \in [-\lambda, \lambda]$$

Proximity Operator: per coordinate

$$\left[\text{prox}_{\gamma \lambda \|\cdot\|_1}(u) \right]_i = \begin{cases} u_i - \lambda\gamma & \text{if } u_i > \lambda\gamma \\ 0 & \text{if } u_i \in [-\lambda\gamma, \lambda\gamma] \\ u_i + \lambda\gamma & \text{if } u_i < -\lambda\gamma \end{cases}$$

Proximal Gradient (aka ISTA):

$$\begin{cases} u_{k+1} = x_k - \gamma A^\top (Ax_k - b) \\ x_{k+1} = \text{prox}_{\gamma \lambda \|\cdot\|_1}(u_{k+1}) \end{cases}$$



Example: LASSO

$$\begin{aligned} \text{Find } x^* \in \arg \min_{x \in \mathbb{R}^n} \mathcal{R}(x; \{a_i, b_i\}_{i=1}^m) &+ \lambda r(x) \\ \text{Find } x^* \in \arg \min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2 &+ \lambda \|x\|_1 \end{aligned}$$

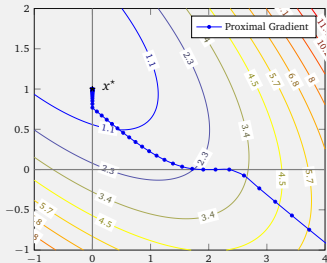
smooth non-smooth

Coordinates	Structure	\leftrightarrow	Optimality conditions	\leftrightarrow	Proximity operation
$\forall i$	$x_i^* = 0$	\Leftrightarrow	$A_i^\top (Ax^* - b) \in [-\lambda, \lambda]$	\Leftrightarrow	$\left[\text{prox}_{\gamma \lambda \ \cdot\ _1}(u^*) \right]_i = 0$ $u^* = x^* - \gamma A^\top (Ax^* - b)$

$$\left[\text{prox}_{\gamma \lambda \|\cdot\|_1}(u) \right]_i = \begin{cases} u_i - \lambda\gamma & \text{if } u_i > \lambda\gamma \\ 0 & \text{if } u_i \in [-\lambda\gamma; \lambda\gamma] \\ u_i + \lambda\gamma & \text{if } u_i < -\lambda\gamma \end{cases}$$

Proximal Gradient (aka ISTA):

$$\begin{cases} u_{k+1} = x_k - \gamma A^\top (Ax_k - b) \\ x_{k+1} = \text{prox}_{\gamma \lambda \|\cdot\|_1}(u_{k+1}) \end{cases}$$



Iterates (x_k) reach the **same structure as x^*** in **finite** time!

In practice: After that **identification**: fixed structure and better rate

>>> Distributed Proximal Gradient

Algorithm

Worker i updates w/ local data (f_i)

$$x_i^{k+1/2} = x^k - \gamma \nabla f_i(x^k)$$

for all $i = 1, \dots, M$

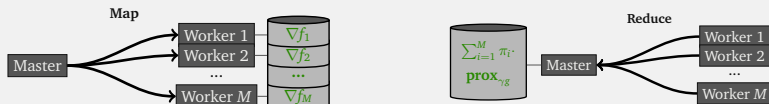
Master gathers the local variables

$$\bar{x}^{k+1} = \sum_{i=1}^M \pi_i x_i^{k+1/2}$$

Master performs a proximity operation

$$x_1^{k+1} = \dots = x_M^{k+1} = \mathbf{prox}_{\gamma g}(\bar{x}^{k+1})$$

Implementation



Distributed Proximal Gradient

Master:

Initialize $\bar{x} = \bar{x}^0$,

while not converged **do**

when all workers have finished:

 Receive (x_i) from each of them

$$\bar{x} \leftarrow \sum_{i=1}^M \pi_i x_i$$

$$x \leftarrow \mathbf{prox}_{\gamma g}(\bar{x})$$

 Broadcast x to all agents

$$k \leftarrow k + 1$$

Interrupt all slaves

Output x

Worker i :

Initialize $x = x_i = \bar{x}$,

while not interrupted by master **do**

 Receive the most recent x

$$x_i \leftarrow x - \gamma \nabla f_i(x)$$

 Send x_i to the master

$$f_i(x) = \frac{1}{|S_i|} \sum_{j \in S_i} \ell_j(x)$$

Local risk at worker i

Communications may soon become the bottleneck in distributed learning, hence the rise of asynchronous methods.

DAve-PG

Master:

```

Initialize  $\bar{x}$ 
while not converged do
  when a worker finishes:
    Receive adjustment  $\Delta$  from it
     $\bar{x} \leftarrow \bar{x} + \Delta$ 
     $x \leftarrow \text{prox}_{\gamma g}(\bar{x})$ 
    Send  $x$  to the agent in return
     $k \leftarrow k + 1$ 
Interrupt all slaves
Output  $x$ 
    
```

Worker i :

```

Initialize  $x = x_i = \bar{x}$ 
while not interrupted by master do
  Receive the most recent  $x$ 
   $x_i \leftarrow x - \gamma \nabla f_i(x)$ 
   $\Delta \leftarrow \pi_i(x_i - x_i^{prev})$     $x_i^{prev} \leftarrow x_i$ 
  Send adjustment  $\Delta$  to master
    
```

$$f_i(x) = \frac{1}{|\mathcal{S}_i|} \sum_{j \in \mathcal{S}_i} \ell_j(x)$$

Local risk at worker i

- > With *sparsity inducing* regularizers (eg. ℓ_1 norm), *master-to-worker* communications will eventually become sparse \Rightarrow **Use it!**
identification of proximal methods
- > Unfortunately, *worker-to-master* communications stay dense...
Idea: sparsify adaptively!

- > Machine Learning problems often have a *noticeable structure*;
sparsity, low rank
- > This structure is identified progressively by *proximal methods*;
+ CD, Var. Red., Distributed methods, etc.
- > For most problem, we *do not know* if the identified structure is optimal;
adaptivity is key
- > Whenever structure appears, it can often be used numerically
storage, communications
- > Importance of conditionning and adaptation frequency
to achieve best theoretical performance

▷ Grishchenko, I., Malick, Amini: *Distributed Learning with Sparse Communications by Identification*, 2020 <https://arxiv.org/abs/1812.03871>

Thanks to IDEX UGA IRS DOLL



& PGMO



Thank you! – Franck IUTZELER <http://www.iutzeler.org>