

# INEXACT SMOOTH PENALTY FOR CONVEX PROBLEMS WITH LINEAR CONSTRAINTS

**Tatiana Tatarenko and Angelia Nedić**

TU Darmstadt  
Arizona State University at Tempe

**November 15, 2021**

## Problem of Interest

- We want to *efficiently* solve the following convex problem

$$\text{minimize } f(x) \quad \text{subject to } \langle a_i, x \rangle \leq b_i, \quad i = 1, \dots, m, \quad (1)$$

- The function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is strongly convex and has Lipschitz continuous gradients
- The number  $m \geq 2$  of inequalities is assumed to be large.
- Such a problem arises in *machine learning* (training a classifier):
  - Each data point induces a constraint
  - $f(\cdot)$  corresponds to a "regularizer" that we choose (such as  $\|\cdot\|^2$ )
- When  $f \equiv 0$ , the randomized Kaczmarz algorithm can solve the resulting feasibility problem with a geometric rate [Strohmer and Vershynin 2008]
- Geometric rate results also exist for a general convex sets  $X_i$  admitting efficient projections [Nedić 2010, Richtarik and Necoara 2018]
- If max-type penalty is used, the resulting penalized problem is non-differentiable:

$$\text{minimize } f(x) + \frac{\gamma}{m} \sum_{i=1}^m \max\{0, \langle a_i, x \rangle - b_i\}$$

- At best, the number  $k$  of iterations is in the order of  $O(1/k)$  (subgradient method)
- Goal: to apply algorithm converging with a geometric rate, i.e. SAGA for example

## Proposed Approach: Basic Idea

- Consider a penalized problem of the form

$$\text{minimize } f(x) + \frac{\gamma}{m} \sum_{i=1}^m h_{\delta}(x; a_i, b_i),$$

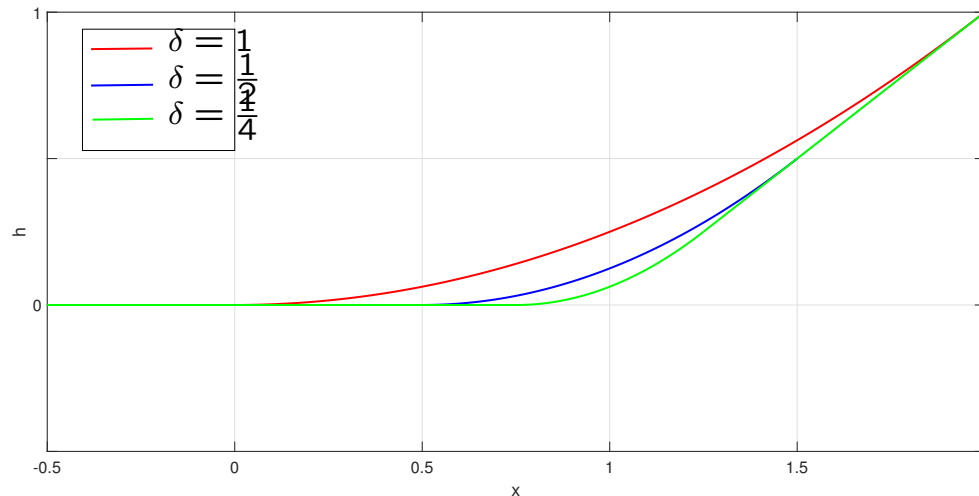
where  $h_{\delta}(x; a_i, b_i)$  is a *convex* penalty function associated with constraint  $\langle a_i, x \rangle \leq b_i$  and  $\gamma > 0$  and  $\delta > 0$  are penalty parameters

- Our choice of the penalty function is guided by the desire that **the resulting penalized objective is strongly convex and has Lipschitz gradients**
- So we *will choose*  $h_{\delta}(x; a_i, b_i)$  to have *Lipschitz continuous gradients*
- Hence, SAGA or other fast incremental approaches can be applied for solving the penalized problem; however, the penalized problem is not necessarily equivalent to the original problem
- **Main goal:** Determine the parameters  $\delta$  and  $\gamma$ , so that **the solution of the penalized problem is**
  - Feasible for the original problem and
  - Sub-optimal for the original problem with a desired pre-specified accuracy
- It turns out, we can accomplish this by choosing:
  - *Inexact penalty function that induces a positive penalty inside the feasible set*
- **The approach will work for the class of functions that have bounded level sets**

## Penalized Problem

$$\min_x F_{\gamma\delta}(x), \quad F_{\gamma\delta}(x) = f(x) + \frac{\gamma}{m} \sum_{i=1}^m h_{\delta}(x; a_i, b_i),$$

$$h_{\delta}(x; a, b) = \begin{cases} \frac{\langle a, x \rangle - b}{\|a\|}, & \text{if } \langle a, x \rangle - b > \delta, \\ \frac{(\langle a, x \rangle - b + \delta)^2}{4\delta\|a\|}, & \text{if } -\delta \leq \langle a, x \rangle - b \leq \delta, \\ 0, & \text{if } \langle a, x \rangle - b < -\delta. \end{cases}$$



**Figure 1:** Penalty function  $h_{\delta}(x; 1, 1)$  for the constraint  $x \leq 1$ ,  $x \in \mathbb{R}$ . The case  $\delta = 0$  corresponds to  $h_0(x; 1, 1) = \max\{0, x - 1\}$

Some properties of the penalty function  $h_\delta(x; a, b)$ :

- $h_0(x; a, b) = \text{dist}(x, X_i)$ ,  $X_i = \{x \mid \langle a_i, x \rangle \leq b_i\}$ ,  $i = 1, \dots, m$ .
- If  $0 < \delta \leq \delta'$ , then  $h_\delta(x; a, b) \leq h_{\delta'}(x; a, b)$  for all  $x \in \mathbb{R}^n$ .
- For  $\delta > 0$ , we have  $\|\nabla h_\delta(x; a, b) - \nabla h_\delta(y; a, b)\| \leq \frac{1}{2\delta} \|x - y\|$  for all  $x, y \in \mathbb{R}^n$ .
- If  $x$  satisfies  $\langle a_i, x \rangle \leq b_i$ , then  $h_\delta(x; a_i, b_i) \leq \frac{\delta}{4\|a_i\|}$ .

An important consequence: nested level sets for any  $\gamma > 0$ ,  $\delta < \delta'$ , and any  $t \in \mathbb{R}$ ,

$$\{x \mid F_{\gamma\delta'}(x) \leq t\} \subseteq \{x \mid F_{\gamma\delta}(x) \leq t\} \subseteq \{x \mid f(x) \leq t\}.$$

Thus, if  $f$  has bounded level sets, so does  $F_{\gamma\delta}(x)$  for any  $\gamma > 0$  and  $\delta > 0$ .

### Assumption

- The feasible set  $X = \{x \mid \langle a_i, x \rangle \leq b_i, i = 1, \dots, m\}$  has a nonempty interior:  
there is  $\hat{x}$  such that:  $\langle a_i, \hat{x} \rangle \leq b_i - \epsilon$  for all  $i = 1, \dots, m$ .

We let  $X_i = \{x \mid \langle a_i, x \rangle \leq b_i\}$ , so that  $X = \bigcap_{i=1}^m X_i$ .

Given a system of linear equations/inequalities, there is  $\beta > 0$  such that

$$\beta \sum_{i=1}^m \text{dist}(x, X_i) \geq \text{dist}(x, X) \quad \text{for all } x \in \mathbb{R}^n,$$

$\beta$  is a **Hoffman bound** (O. Guler, A. Hoffman, and U. Rothblum 1992), and  $\text{dist}(x, Y) = \|x - \Pi_Y[x]\|$  for a convex closed set  $Y \subset \mathbb{R}^n$  and any  $x \in \mathbb{R}^n$

For our affine set  $X$ , Hoffman bound  $\beta$  depends on the vectors  $a_i$

## Strongly Convex $f$

- Let  $x^*$  be the optimal solution of the original problem, i.e.,  $x^* = \operatorname{argmin}_{x \in X} f(x)$
- Let  $\epsilon_a$  be the desired accuracy for solving the problem

**Proposition 1** *Let  $X$  have nonempty interior, and let  $f$  be strongly convex with a constant  $\mu > 0$ . Let  $\delta$  and  $\gamma$  be such that*

$$0 < \delta < \min \left\{ \epsilon, \frac{16\alpha_{\min}^2}{\beta^2 m^2} \right\}, \quad \gamma \leq \frac{2\mu}{\delta} \epsilon_a \alpha_{\min}$$

$$\gamma \geq \max \left\{ L \left( \frac{1}{m\beta} - \frac{\sqrt{\delta}}{4\alpha_{\min}} \right)^{-1}, 4mL\alpha_{\max} \left( \frac{1}{\sqrt{\delta}} + \frac{\beta m}{\alpha_{\min}} \right) \right\},$$

where  $L$  is a Lipschitz constant for  $f$  over some suitably defined level set,  $\alpha_{\min} = \min_i \|a_i\|$  and  $\alpha_{\max} = \max_i \|a_i\|$ . Then, the solution  $x_{\gamma\delta}^*$  of the penalized problem  $\min_x F_{\gamma\delta}(x)$  satisfies

$$x_{\gamma\delta}^* \in X, \quad \|x_{\gamma\delta}^* - x^*\|^2 \leq \epsilon_a.$$

## Proof: Main Steps

- When the vector  $x_{\gamma\delta}^*$  is feasible i.e.,  $x_{\gamma\delta}^* \in X$ , we have

$$f(x^*) \leq f(x_{\gamma\delta}^*). \quad (2)$$

Since the penalty functions are non-negative, we have  $h_\delta(x_{\gamma\delta}^*; a_i, b_i) \geq 0$  for all  $i = 1, \dots, m$ . The point  $x^*$  is feasible but it may be penalized, in which case  $h_\delta(x^*; a_i, b_i) \leq \frac{\delta}{4\|a_i\|}$ . Therefore, we have

$$h_\delta(x^*; a_i, b_i) - h_\delta(x_{\gamma\delta}^*; a_i, b_i) \leq \frac{\delta}{4\|a_i\|} \quad \text{for all } i = 1, \dots, m. \quad (3)$$

Using the relations (2) and (3), we obtain

$$\begin{aligned} F_{\gamma\delta}(x^*) - F_{\gamma\delta}(x_{\gamma\delta}^*) &= f(x^*) - f(x_{\gamma\delta}^*) + \frac{\gamma}{m} \left( \sum_{i=1}^m h_\delta(x^*; a_i, b_i) - h_\delta(x_{\gamma\delta}^*; a_i, b_i) \right) \\ &\leq \frac{\gamma\delta}{4\alpha_{\min}}. \end{aligned}$$

By the strong convexity of  $F_{\gamma\delta}$  and the fact that  $x_{\gamma\delta}^*$  is the minimum of  $F_{\gamma\delta}(x)$ , it follows that

$$\|x^* - x_{\gamma\delta}^*\|^2 \leq \frac{2}{\mu_f} (F_{\gamma\delta}(x^*) - F_{\gamma\delta}(x_{\gamma\delta}^*)) \leq \frac{\gamma\delta}{2\mu_f\alpha_{\min}} \leq \epsilon_a, \quad (4)$$

where the last inequality in the preceding relation is due to the choice of  $\gamma \leq \frac{2\mu f}{\delta} \alpha_{\min} \epsilon_a$ .

- Proving the feasibility of  $x_{\gamma\delta}^*$  is much more involved. Define

$$\hat{x}_{\gamma\delta}^* = \Pi_X[x_{\gamma\delta}^*].$$

We consider two possibilities:  $\|\hat{x}_{\gamma\delta}^* - x_{\gamma\delta}^*\| \geq \sqrt{\delta}$  and  $\|\hat{x}_{\gamma\delta}^* - x_{\gamma\delta}^*\| < \sqrt{\delta}$ .

*Case 1:*  $\|\hat{x}_{\gamma\delta}^* - x_{\gamma\delta}^*\| > \sqrt{\delta}$ . Recall that we have  $h_\delta(x; a_i, b_i) \geq \text{dist}(x, X_i)$  for all  $i$ .

Thus, by the definition of the functions  $F_{\gamma\delta}$ , for any  $x \in \mathbb{R}^n$  we can write

$$F_{\gamma\delta}(x) \geq f(x) + \frac{\gamma}{m} \sum_{i=1}^m \text{dist}(x, X_i).$$

Then, by Hoffman's lemma, for some  $\beta > 0$  we have

$$F_{\gamma\delta}(x) \geq f(x) + \frac{\gamma}{m\beta} \text{dist}(x, X) \quad \text{for all } x \in \mathbb{R}^n.$$

Letting  $x = x_{\gamma\delta}^*$  in the preceding relation, we obtain

$$\begin{aligned} F_{\gamma\delta}(x_{\gamma\delta}^*) &\geq f(x_{\gamma\delta}^*) + \frac{\gamma}{m\beta} \|\hat{x}_{\gamma\delta}^* - x_{\gamma\delta}^*\| + f(\hat{x}_{\gamma\delta}^*) - f(x_{\gamma\delta}^*) \\ &\geq \frac{\gamma}{m\beta} \|\hat{x}_{\gamma\delta}^* - x_{\gamma\delta}^*\| - L \|\hat{x}_{\gamma\delta}^* - x_{\gamma\delta}^*\| + f(\hat{x}_{\gamma\delta}^*) \\ &= \left( \frac{\gamma}{m\beta} - L \right) \|\hat{x}_{\gamma\delta}^* - x_{\gamma\delta}^*\| + F_{\gamma\delta}(\hat{x}_{\gamma\delta}^*) - \frac{\gamma}{m} \sum_{i=1}^m h_\delta(\hat{x}_{\gamma\delta}^*; a_i, b_i), \end{aligned}$$



where in the second inequality we use the assumption that the norms of the subgradients in the subdifferential set  $\partial f(x)$  are bounded by  $L$  in a region containing the point  $x = \hat{x}_{\gamma\delta}^*$  (requires a lemma to assert this). Taking into the account that  $h_\delta(x; a_i, b_i) \leq \frac{\delta}{4\|a_i\|}$  when  $x \in X_i$ , and using  $\hat{x}_{\gamma\delta}^* \in X \subseteq X_i$ , we see that

$$F_{\gamma\delta}(x_{\gamma\delta}^*) \geq \left( \frac{\gamma}{m\beta} - L \right) \|x_{\gamma\delta}^* - \hat{x}_{\gamma\delta}^*\| + F_{\gamma\delta}(\hat{x}_{\gamma\delta}^*) - \frac{\gamma\delta}{4m} \sum_{i=1}^m \frac{1}{\|a_i\|}.$$

The conditions on  $\gamma$  imply that  $\gamma \geq Lm\beta$ . Using the relations  $\gamma \geq Lm\beta$  and  $\|\hat{x}_{\gamma\delta}^* - x_{\gamma\delta}^*\| > \sqrt{\delta}$ , we further obtain

$$F_{\gamma\delta}(x_{\gamma\delta}^*) > \left( \frac{\gamma}{m\beta} - L \right) \sqrt{\delta} - \frac{\gamma\delta}{4\alpha_{\min}} + F_{\gamma\delta}(\hat{x}_{\gamma\delta}^*) \geq F_{\gamma\delta}(\hat{x}_{\gamma\delta}^*), \quad (5)$$

where the last inequality is obtained by using  $\left( \frac{\gamma}{m\beta} - L \right) \sqrt{\delta} - \frac{\gamma\delta}{4\alpha_{\min}} \geq 0$ , equivalent

$$\frac{\gamma}{m\beta} - L - \frac{\gamma\sqrt{\delta}}{4\alpha_{\min}} \geq 0 \quad \iff \quad \gamma \left( \frac{1}{m\beta} - \frac{\sqrt{\delta}}{4\alpha_{\min}} \right) \geq L.$$

The last inequality holds by the conditions imposed on the parameters  $\gamma$  and  $\delta$ . Thus, relation (5) implies that

$$F_{\gamma\delta}(x_{\gamma\delta}^*) > F_{\gamma\delta}(\hat{x}_{\gamma\delta}^*),$$

which contradicts the fact that  $x_{\gamma\delta}^*$  is an unconstrained minimizer of  $F_{\gamma\delta}$ .

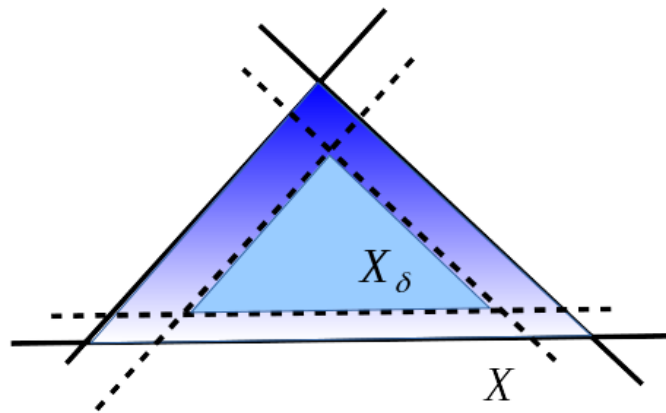
Case 2:  $\|\hat{x}_{\gamma\delta}^* - x_{\gamma\delta}^*\| \leq \sqrt{\delta}$ . Uses the interior-point assumption and a special construction. Specifically, it relies on the following result.

**Lemma 1** *Let the interior-point assumption hold and let  $\delta$  be such that  $0 < \delta \leq \epsilon$ . Then, for any  $x \notin X$  there exists a feasible point  $x_{in} \in X$  such that*

$$(a) \ h_{\delta}(x_{in}; a_j, b_j) = 0 \quad \text{for all } j = 1, \dots, m,$$

$$(b) \ \|x - x_{in}\| \leq \|x - \Pi_X[x]\| + \frac{\beta m \delta}{\min_{1 \leq i \leq m} \|a_i\|},$$

where  $\beta$  is Hoffman's bound.



**Figure 2:** Illustration of the set  $X_\delta$ .

## Connecting with SAGA by Defazio et al. 2014

---

### Algorithm 1 SAGA-based Fast Incremental Method for Solving Penalized Problem

---

- Let  $x^0 \in \mathbb{R}^n$  and  $\nabla f(\phi_i^0) + \gamma \nabla h_\delta(\phi_i^0; a_i, b_i)$  with  $\phi_i^0 = x^0$  be available for  $i = 1, \dots, m$ .
  - Pick an index  $j \in \{1, 2, \dots, m\}$  uniformly at random.
  - Set  $\phi_j^{t+1} = x^t$  and store  $\nabla f(\phi_j^{t+1}) + \gamma \nabla h_\delta(\phi_j^{t+1}; a_j, b_j)$ .
  - $x^{t+1} = x^t - \alpha [\nabla f(\phi_j^{t+1}) + \gamma \nabla h_\delta(\phi_j^{t+1}; a_j, b_j) - \nabla f(\phi_j^t) - \gamma \nabla h_\delta(\phi_j^t; a_j, b_j) + \frac{1}{m} \sum_{i=1}^m (\nabla f(\phi_i^t) + \gamma \nabla h_\delta(\phi_i^t; a_i, b_i))]$ .
-

## Applying SAGA to the Penalized Problem

**Proposition 2** *Let the interior-point assumption hold, and let the function  $f$  be strongly convex with a parameter  $\mu > 0$  and have Lipschitz continuous gradients with a constant  $L_f > 0$ . Let  $x^* = \operatorname{argmin}_{x \in X} f(x)$ , and assume that an accuracy level  $\epsilon_a$  is given. Consider Algorithm SAGA applied to the penalized problem:*

$$\min_x \frac{1}{m} \sum_{i=1}^m (f(x) + \gamma h_\delta(x; a_i, b_i)),$$

where the penalty parameters  $\gamma$  and  $\delta$  are chosen to satisfy the conditions of Proposition 1, and the step size is given by

$$\alpha = \frac{1}{2(\mu m + L_f + \frac{\gamma \alpha_{\max}}{2\delta})},$$

with  $\alpha_{\max} = \max_i \|a_i\|$ .

Then, the following convergence rate result is valid for the iterates of the algorithm:

$$\mathbb{E}[\|x^t - \Pi_X[x^t]\|^2] \leq O(q_\gamma^t), \quad \mathbb{E}[\|x^t - x^*\|^2] \leq O(q_\gamma^t) + 2\epsilon_a,$$

$$q_\gamma = 1 - \frac{\mu}{2(\mu m + L_f + \frac{\gamma \alpha_{\max}}{2\delta})}.$$

Proof: Immediately from Proposition 1 and a rate result from Defazio et al. 2014.

## Convergence Result for Merely Convex Function $f$

**Proposition 3** *Let the interior-point assumption hold. Let  $f$  be convex function with bounded level sets and Lipschitz continuous gradients with a constant  $L_f$ , and let  $f^* = \min_{x \in X} f(x)$ . Let  $\epsilon_a$  a desired accuracy. Consider the iterates  $x_k$  produced by SAGA method applied to the penalized problem:*

$$\min_x \frac{1}{m} \sum_{i=1}^m (f(x) + \gamma h_\delta(x; a_i, b_i)),$$

where  $\gamma$  and  $\delta$  are chosen such that the conditions of Proposition 1 are satisfied, and the stepsize is given by

$$\alpha = \frac{1}{3(L_f + \gamma\alpha_{\max}/(2\delta))}$$

Then, for the iterate averages

$$\bar{x}_k = \frac{1}{k} \sum_{i=0}^{k-1} x_i$$

the following hold relation holds

$$\mathbb{E}[f(\bar{x}_k)] - f^* \leq O(1/k) + 2\epsilon_a \quad \text{for all } k,$$

Moreover, there exists some  $T > 0$  such that for all  $k$ ,

$$-\gamma\epsilon_a \leq \mathbb{E}[f(\bar{x}_k)] - f^*$$

## Simulation Results

$$f(x) = \frac{1}{2} \|\Phi x - x^0\|^2 \quad Ax \leq b, \quad A \in \mathbb{R}^{m \times n}, \quad n = 30$$

where  $x^0$  is chosen randomly according to a zero-mean normal distribution with covariance matrix  $\Sigma = 10I$ . Constraint set  $X$  is constructed so that the optimal point  $x^0$  is on its boundary. The SAGA algorithm applied to the penalized problem is compared to an algorithm that, at time  $t$ , selects one constraint  $X_{\omega_t}$  randomly from the collection  $X_1, \dots, X_m$ , and performs a gradient update\*

$$x^{t+1} = \Pi_{X_{\omega_t}}[x^t - \alpha_t \nabla f(x^t)],$$

where  $\alpha_t$  is a diminishing stepsize (not summable, but square summable).

The full gradient method is also simulated.

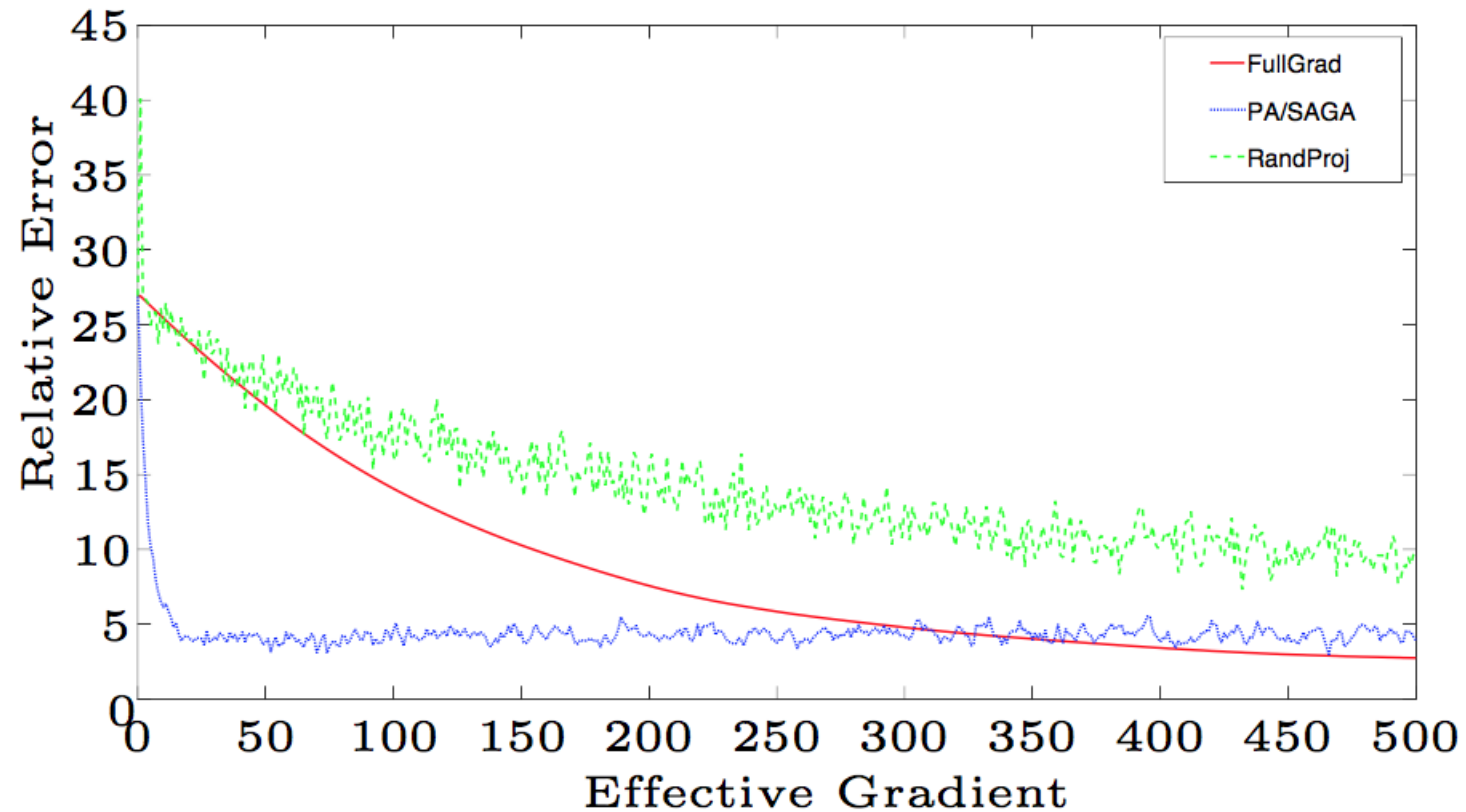
The algorithms are run for 1,000 iterations, and the relative error is plotted

$$\frac{\|x^t - x^*\|}{\|x^*\|}$$

along the iterate sequence.

---

\*Nedić 2011



**Figure 3: Full-Gradient, PA/SAGA (our penalized approach) and Random Projections,  $m = 1500$ .**

## Summary and Extensions

- We provided a novel penalty re-formulation for a convex minimization problem with linear constraints.
- The structure of the penalty functions that we used to penalize the linear constraints, and the suitable choices of the penalty parameters render the penalized unconstrained problem with solutions that are *feasible* for the original constrained problem.
- With an additional constraint on the penalty parameters imposed by a desired accuracy level, the solutions of the penalized unconstrained problem are guaranteed to be arbitrarily close to the solution set of the original problem.
- An advantage of the proposed penalty reformulation is in the ability to employ fast incremental gradient methods, such as SAGA.
- We have a convergence result for our penalty approach also when  $f$  has bounded level sets (not necessarily strongly convex)



- Difficulty: our results rely on the availability of Hoffman constant  $\beta$ ,

$$\beta \sum_{i=1}^m \text{dist}(x, X_i) \geq \text{dist}(x, X) \quad \text{for all } x \in \mathbb{R}^n,$$

which is hard to (upper) estimate; see recent work by J.Pena, J. Vera L.F. Zuluaga *New characterizations of Hoffman constants for systems of linear constraints*, *Mathematical Programming* 187, p. 79–109, 2021, arxiv <https://arxiv.org/pdf/1905.02894.pdf>

- As a remedy, we have considered an approach where we vary the penalty parameter  $\gamma$  with time at each iteration and show that such a method converges to the optimal point

- Details are in our paper:

T. Tatarenko and AN *A Smooth Inexact Penalty Reformulation of Convex Problems with Linear Constraints*, *SIAM J. on Optim.*, 31 (3), 2141–2170, 2021