

# Accelerated and Sparse Algorithms for Approximate Personalized PageRank

**David Martínez-Rubio**

joint work with Elias Wirth, Sebastian Pokutta

Technische Universität Berlin, Zuse Institute Berlin



# Problem

Problem:

$$\min_{x \in \mathbb{R}_{\geq 0}^n} \{g(x) \stackrel{\text{def}}{=} \langle x, Qx \rangle - \langle b, x \rangle\}.$$

where  $0 \prec \alpha \cdot I \preccurlyeq Q \preccurlyeq L \cdot I$  and  $Q_{ij} \leq 0$ .

# Problem

Problem:

$$\min_{\mathbf{x} \in \mathbb{R}_{\geq 0}^n} \{g(\mathbf{x}) \stackrel{\text{def}}{=} \langle \mathbf{x}, Q\mathbf{x} \rangle - \langle \mathbf{b}, \mathbf{x} \rangle\}.$$

where  $0 \prec \alpha \cdot I \preceq Q \preceq L \cdot I$  and  $Q_{ij} \leq 0$ . For  $\ell_1$ -regularized personalized PageRank, it is

$$Q \stackrel{\text{def}}{=} \alpha I + \frac{1-\alpha}{2} \mathcal{L} \quad \text{and} \quad \mathbf{b} \stackrel{\text{def}}{=} \alpha \left( D^{-1/2} \mathbf{s} - \rho D^{1/2} \mathbf{1} \right)$$

where  $\alpha, \rho > 0$ ,  $\mathcal{L} \stackrel{\text{def}}{=} I - D^{-1/2} A D^{-1/2}$  is the symmetric normalized Laplacian matrix, which satisfies  $0 \prec \mathcal{L} \preceq 2I$ .

# Problem

Problem:

$$\min_{\mathbf{x} \in \mathbb{R}_{\geq 0}^n} \{g(\mathbf{x}) \stackrel{\text{def}}{=} \langle \mathbf{x}, Q\mathbf{x} \rangle - \langle \mathbf{b}, \mathbf{x} \rangle\}.$$

where  $0 \prec \alpha \cdot I \preceq Q \preceq L \cdot I$  and  $Q_{ij} \leq 0$ . For  $\ell_1$ -regularized personalized PageRank, it is

$$Q \stackrel{\text{def}}{=} \alpha I + \frac{1-\alpha}{2} \mathcal{L} \quad \text{and} \quad \mathbf{b} \stackrel{\text{def}}{=} \alpha \left( D^{-1/2} \mathbf{s} - \rho D^{1/2} \mathbf{1} \right)$$

where  $\alpha, \rho > 0$ ,  $\mathcal{L} \stackrel{\text{def}}{=} I - D^{-1/2} A D^{-1/2}$  is the symmetric normalized Laplacian matrix, which satisfies  $0 \prec \mathcal{L} \preceq 2I$ .

The problem comes from the personalized PageRank problem

$$\min \{f(\mathbf{x}) \stackrel{\text{def}}{=} \langle \mathbf{x}, Q\mathbf{x} \rangle - \alpha \langle D^{-1/2} \mathbf{s}, \mathbf{x} \rangle\},$$

by adding the  $\ell_1$  regularization  $+\alpha\rho\|D^{1/2}\mathbf{x}\|_1$  and noticing that the minimizer is in  $\mathbb{R}_{\geq 0}$ , for  $\rho > 0$ . The personalized PageRank vector is the solution to the system

$$\mathbf{x} = (1 - \alpha)W\mathbf{x} + \alpha\mathbf{s} = ((1 - \alpha)W + \alpha\mathbf{s}\mathbf{1}^T)\mathbf{x},$$

where  $W = (I + AD^{-1})/2$  and  $\mathbf{s} \in \Delta^n$  is a distribution over the nodes.

# Results and comparison

- ▶ The Hessian of  $g$  is  $Q$ , satisfying  $\alpha I \preceq Q \preceq LI$ , its condition number is  $L/\alpha$ .
- ▶  $\mathcal{S}^* \stackrel{\text{def}}{=} \text{supp}(\mathbf{x}^*)$ ,  $\text{vol}(\mathcal{S}^*) \stackrel{\text{def}}{=} \text{nnz}(Q_{:, \mathcal{S}^*})$  and  $\widetilde{\text{vol}}(\mathcal{S}^*) \stackrel{\text{def}}{=} \text{nnz}(Q_{\mathcal{S}^*, \mathcal{S}^*})$ .
- ▶ For the  $\ell_1$ -regularized personalized PageRank, it is  $\text{vol}(\mathcal{S}^*) \leq \frac{1}{\rho} + |\mathcal{S}^*|$ .

Method	Time complexity	Space complexity
ISTA [FRS+19]	$\widetilde{\mathcal{O}}(\text{vol}(\mathcal{S}^*) \frac{L}{\alpha})$	$\mathcal{O}( \mathcal{S}^* )$
CDPR ( <b>Ours</b> )	$\mathcal{O}( \mathcal{S}^* ^3 +  \mathcal{S}^*  \text{vol}(\mathcal{S}^*))$	$\mathcal{O}( \mathcal{S}^* ^2)$
ASPR ( <b>Ours</b> )	$\widetilde{\mathcal{O}}( \mathcal{S}^*  \widetilde{\text{vol}}(\mathcal{S}^*) \sqrt{\frac{L}{\alpha}} +  \mathcal{S}^*  \text{vol}(\mathcal{S}^*))$	$\mathcal{O}( \mathcal{S}^* )$
CASPR ( <b>Ours</b> )	$\widetilde{\mathcal{O}}( \mathcal{S}^*  \widetilde{\text{vol}}(\mathcal{S}^*) \min \left\{ \sqrt{\frac{L}{\alpha}},  \mathcal{S}^*  \right\} +  \mathcal{S}^*  \text{vol}(\mathcal{S}^*))$	$\mathcal{O}( \mathcal{S}^* )$

- ▶ Problem:

$$\min_{\mathbf{x} \in \mathbb{R}_{\geq \mathbf{0}}^n} \{g(\mathbf{x}) \stackrel{\text{def}}{=} \langle \mathbf{x}, Q\mathbf{x} \rangle - \langle \mathbf{b}, \mathbf{x} \rangle\}.$$

# A geometric lemma

Suppose:

- ▶  $x^{(0)} \in \mathbb{R}_{\geq 0}^n$  and  $S \subseteq [n]$  s.t.  $x_i^{(0)} = 0$  if  $i \notin S$  and  $\nabla_i g(x^{(0)}) \leq 0$  if  $i \in S$ .
- ▶  $C \stackrel{\text{def}}{=} \text{span}(\{e_i \mid i \in S\}) \cap \mathbb{R}_{\geq 0}^n$ .
- ▶  $x^{(*,C)} \stackrel{\text{def}}{=} \arg \min_{x \in C} g(x)$  and  $x^* \stackrel{\text{def}}{=} \arg \min_{x \in \mathbb{R}_{\geq 0}^n} g(x)$ .

# A geometric lemma

Suppose:

- ▶  $\mathbf{x}^{(0)} \in \mathbb{R}_{\geq 0}^n$  and  $S \subseteq [n]$  s.t.  $x_i^{(0)} = 0$  if  $i \notin S$  and  $\nabla_i g(\mathbf{x}^{(0)}) \leq 0$  if  $i \in S$ .
- ▶  $C \stackrel{\text{def}}{=} \text{span}(\{\mathbf{e}_i \mid i \in S\}) \cap \mathbb{R}_{\geq 0}^n$ .
- ▶  $\mathbf{x}^{(*,C)} \stackrel{\text{def}}{=} \arg \min_{\mathbf{x} \in C} g(\mathbf{x})$  and  $\mathbf{x}^* \stackrel{\text{def}}{=} \arg \min_{\mathbf{x} \in \mathbb{R}_{\geq 0}^n} g(\mathbf{x})$ .

Then:

1. It holds that  $\mathbf{x}^{(0)} \leq \mathbf{x}^{(*,C)}$  and  $\nabla_i g(\mathbf{x}^{(*,C)}) = 0$  for all  $i \in S$ .
2. If for  $i \in S$ , we have  $x_i^{(0)} > 0$  or  $\nabla_i g(\mathbf{x}^{(0)}) < 0$ , then  $x_i^{(*,C)} > 0$ .
3. If  $x_i^{(*,C)} > 0$  for all  $i \in S$ , we have  $\mathbf{x}^{(*,C)} \leq \mathbf{x}^*$  and therefore  $S \subseteq S^*$ .

# A geometric lemma

Suppose:

- ▶  $x^{(0)} \in \mathbb{R}_{\geq 0}^n$  and  $S \subseteq [n]$  s.t.  $x_i^{(0)} = 0$  if  $i \notin S$  and  $\nabla_i g(x^{(0)}) \leq 0$  if  $i \in S$ .
- ▶  $C \stackrel{\text{def}}{=} \text{span}(\{e_i \mid i \in S\}) \cap \mathbb{R}_{\geq 0}^n$ .
- ▶  $x^{(*,C)} \stackrel{\text{def}}{=} \arg \min_{x \in C} g(x)$  and  $x^* \stackrel{\text{def}}{=} \arg \min_{x \in \mathbb{R}_{\geq 0}^n} g(x)$ .

Then:

1. It holds that  $x^{(0)} \leq x^{(*,C)}$  and  $\nabla_i g(x^{(*,C)}) = 0$  for all  $i \in S$ .
2. If for  $i \in S$ , we have  $x_i^{(0)} > 0$  or  $\nabla_i g(x^{(0)}) < 0$ , then  $x_i^{(*,C)} > 0$ .
3. If  $x_i^{(*,C)} > 0$  for all  $i \in S$ , we have  $x^{(*,C)} \leq x^*$  and therefore  $S \subseteq S^*$ .

**Proof of 1.:**  $\bar{g} \stackrel{\text{def}}{=} g$  restricted to  $\text{span}(\{e_i \mid i \in S\})$ . Let  $\{x^{(t)}\}_{t=0}^\infty$  be the iterates of  $\text{PGD}(C, x^{(0)}, \bar{g})$ . We start with  $\nabla \bar{g}(x^{(0)}) \leq 0$ . By induction:

$$x^{(t+1)} = x^{(t)} - \frac{1}{L} \nabla \bar{g}(x^{(t)}) \geq x^{(t)} \text{ and } \nabla \bar{g}(x^{(t+1)}) = \nabla \bar{g}(x^{(t)}) - \frac{1}{L} Q_{S,S} \nabla \bar{g}(x^{(t)}) \leq 0$$

$x^{(t)} \rightarrow x^{(*,C)}$ ,  $\nabla \bar{g}(x^{(t)}) \rightarrow \nabla \bar{g}(x^{(*,C)})$  (so  $\leq 0$ , and by optimality it is  $\geq 0$ .)



# A geometric lemma

Suppose:

- ▶  $\mathbf{x}^{(0)} \in \mathbb{R}_{\geq 0}^n$  and  $S \subseteq [n]$  s.t.  $x_i^{(0)} = 0$  if  $i \notin S$  and  $\nabla_i g(\mathbf{x}^{(0)}) \leq 0$  if  $i \in S$ .
- ▶  $C \stackrel{\text{def}}{=} \text{span}(\{\mathbf{e}_i \mid i \in S\}) \cap \mathbb{R}_{\geq 0}^n$ .
- ▶  $\mathbf{x}^{(*,C)} \stackrel{\text{def}}{=} \arg \min_{\mathbf{x} \in C} g(\mathbf{x})$  and  $\mathbf{x}^* \stackrel{\text{def}}{=} \arg \min_{\mathbf{x} \in \mathbb{R}_{\geq 0}^n} g(\mathbf{x})$ .

Then:

1. It holds that  $\mathbf{x}^{(0)} \leq \mathbf{x}^{(*,C)}$  and  $\nabla_i g(\mathbf{x}^{(*,C)}) = 0$  for all  $i \in S$ .
2. If for  $i \in S$ , we have  $x_i^{(0)} > 0$  or  $\nabla_i g(\mathbf{x}^{(0)}) < 0$ , then  $x_i^{(*,C)} > 0$ .
3. If  $x_i^{(*,C)} > 0$  for all  $i \in S$ , we have  $\mathbf{x}^{(*,C)} \leq \mathbf{x}^*$  and therefore  $S \subseteq S^*$ .

**Proof of 2.:** We have that  $x_i^{(1)} > 0$  by the assumption on  $x_i^{(0)}$  and the PGD update rule. By the monotonicity of iterates in the proof of 1., we obtain the result.

**Proof of 3.:** Sketch: Apply 1. and 2. to the initial point  $\mathbf{x}^{(*,C)}$  and set of indices  $S \cup \{i \mid \nabla_i g(\mathbf{x}^{(*,C)}) < 0\}$  and then again and so on until you get to  $\mathbf{x}^*$ .

- **Definition.**  $i$  is a good coordinate iff  $i \in \mathcal{S}^*$ . Otherwise it is bad.

# Algorithmic scheme

- ▶ **Definition.**  $i$  is a good coordinate iff  $i \in \mathcal{S}^*$ . Otherwise it is bad.
- ▶ **Idea for an algorithm:** discover good coordinates sequentially, by optimizing in the subspace  $\mathcal{C}^{(t)} \stackrel{\text{def}}{=} \text{span}(\{e_i \mid i \in \mathcal{S}^{(t)}\}) \cap \mathbb{R}_{\geq 0}^n$ , where  $\mathcal{S}^{(t)}$  is the set of currently known good coordinates.

# Algorithmic scheme

- ▶ **Definition.**  $i$  is a good coordinate iff  $i \in \mathcal{S}^*$ . Otherwise it is bad.
- ▶ **Idea for an algorithm:** discover good coordinates sequentially, by optimizing in the subspace  $\mathcal{C}^{(t)} \stackrel{\text{def}}{=} \text{span}(\{e_i \mid i \in \mathcal{S}^{(t)}\}) \cap \mathbb{R}_{\geq 0}^n$ , where  $\mathcal{S}^{(t)}$  is the set of currently known good coordinates.
- ▶ At the minimizer  $x^{(*,t+1)} \stackrel{\text{def}}{=} x^{(*, \mathcal{C}^{(t)})}$ , we are optimal ( $x^{(*,t+1)} = x^*$ ) or we have  $\nabla_i g(x^{(*,t+1)}) < 0$  only if  $i$  is good and new, i.e., only if  $i \in \mathcal{S}^* \setminus \mathcal{S}^{(t)}$ .

# Algorithmic scheme

- ▶ **Definition.**  $i$  is a good coordinate iff  $i \in \mathcal{S}^*$ . Otherwise it is bad.
- ▶ **Idea for an algorithm:** discover good coordinates sequentially, by optimizing in the subspace  $\mathcal{C}^{(t)} \stackrel{\text{def}}{=} \text{span}(\{e_i \mid i \in \mathcal{S}^{(t)}\}) \cap \mathbb{R}_{\geq 0}^n$ , where  $\mathcal{S}^{(t)}$  is the set of currently known good coordinates.
- ▶ At the minimizer  $x^{(*,t+1)} \stackrel{\text{def}}{=} x^{(*, \mathcal{C}^{(t)})}$ , we are optimal ( $x^{(*,t+1)} = x^*$ ) or we have  $\nabla_i g(x^{(*,t+1)}) < 0$  only if  $i$  is good and new, i.e., only if  $i \in \mathcal{S}^* \setminus \mathcal{S}^{(t)}$ .
- ▶ An approximate version of this holds, after overcoming some technicalities.

## An exact algorithm: Conjugate Directions for PageRank (CDPR)

- ▶ Start at  $x^{(0)} = \mathbf{0}$ .

## An exact algorithm: Conjugate Directions for PageRank (CDPR)

- ▶ Start at  $\mathbf{x}^{(0)} = \mathbf{0}$ .
- ▶ For  $t > 0$ , define the set of new good coordinates  $N^{(t)} \stackrel{\text{def}}{=} \{i \in [n] \mid \nabla_i g(\mathbf{x}^{(t)}) < 0\}$  and select  $i \in N^{(t)}$ ,  $\mathbf{u}^{(t)} \stackrel{\text{def}}{=} \nabla_i g(\mathbf{x}^{(t)}) \mathbf{e}_i$ .

## An exact algorithm: Conjugate Directions for PageRank (CDPR)

- ▶ Start at  $\mathbf{x}^{(0)} = \mathbf{0}$ .
- ▶ For  $t > 0$ , define the set of new good coordinates  $N^{(t)} \stackrel{\text{def}}{=} \{i \in [n] \mid \nabla_i g(\mathbf{x}^{(t)}) < 0\}$  and select  $i \in N^{(t)}$ ,  $\mathbf{u}^{(t)} \stackrel{\text{def}}{=} \nabla_i g(\mathbf{x}^{(t)}) \mathbf{e}_i$ .
- ▶ Compute direction  $\mathbf{d}^{(t)}$  from  $\mathbf{u}^{(t)}$  by  $Q$ -Gram-Schmidt using all previous (sparse) directions so  $\langle \mathbf{d}^{(t)}, Q\mathbf{d}^{(k)} \rangle = 0$  for all  $k < t$ .



## An exact algorithm: Conjugate Directions for PageRank (CDPR)

- ▶ Start at  $\mathbf{x}^{(0)} = \mathbf{0}$ .
- ▶ For  $t > 0$ , define the set of new good coordinates  $N^{(t)} \stackrel{\text{def}}{=} \{i \in [n] \mid \nabla_i g(\mathbf{x}^{(t)}) < 0\}$  and select  $i \in N^{(t)}$ ,  $\mathbf{u}^{(t)} \stackrel{\text{def}}{=} \nabla_i g(\mathbf{x}^{(t)}) \mathbf{e}_i$ .
- ▶ Compute direction  $\mathbf{d}^{(t)}$  from  $\mathbf{u}^{(t)}$  by  $Q$ -Gram-Schmidt using all previous (sparse) directions so  $\langle \mathbf{d}^{(t)}, Q\mathbf{d}^{(k)} \rangle = 0$  for all  $k < t$ .
- ▶ Optimize on the line  $\mathbf{x}^{(t+1)} \leftarrow \arg \min_{\eta^{(t)}} \{\mathbf{x}^{(t)} + \eta^{(t)} \mathbf{d}^{(t)}\}$ . It is  $\mathbf{x}^{(t+1)} = \mathbf{x}^{(*, C^{(t)})}$ .

## An exact algorithm: Conjugate Directions for PageRank (CDPR)

- ▶ Start at  $\mathbf{x}^{(0)} = \mathbf{0}$ .
- ▶ For  $t > 0$ , define the set of new good coordinates  $N^{(t)} \stackrel{\text{def}}{=} \{i \in [n] \mid \nabla_i g(\mathbf{x}^{(t)}) < 0\}$  and select  $i \in N^{(t)}$ ,  $\mathbf{u}^{(t)} \stackrel{\text{def}}{=} \nabla_i g(\mathbf{x}^{(t)}) \mathbf{e}_i$ .
- ▶ Compute direction  $\mathbf{d}^{(t)}$  from  $\mathbf{u}^{(t)}$  by  $Q$ -Gram-Schmidt using all previous (sparse) directions so  $\langle \mathbf{d}^{(t)}, Q\mathbf{d}^{(k)} \rangle = 0$  for all  $k < t$ .
- ▶ Optimize on the line  $\mathbf{x}^{(t+1)} \leftarrow \arg \min_{\eta^{(t)}} \{\mathbf{x}^{(t)} + \eta^{(t)} \mathbf{d}^{(t)}\}$ . It is  $\mathbf{x}^{(t+1)} = \mathbf{x}^{(*, C^{(t)})}$ .
- ▶ Time complexity  $\mathcal{O}(|\mathcal{S}^*|^3 + |\mathcal{S}^*| \text{vol}(\mathcal{S}^*))$  and space complexity  $\mathcal{O}(|\mathcal{S}^*|^2)$ .

## An inexact algorithm: Accelerated and Sparse PageRank (ASPR)

1. Because  $Q_{ij} \leq 0$  for  $i \neq j$ , for  $y = x - \Delta e_i$ , we have  $\forall j \neq i$ :  
 $\nabla_j g(y) \geq \nabla_j g(x)$  if  $\Delta > 0$  and  $\nabla_j g(y) \leq \nabla_j g(x)$  otherwise.

## An inexact algorithm: Accelerated and Sparse PageRank (ASPR)

1. Because  $Q_{ij} \leq 0$  for  $i \neq j$ , for  $y = x - \Delta e_i$ , we have  $\forall j \neq i$ :  
 $\nabla_j g(y) \geq \nabla_j g(x)$  if  $\Delta > 0$  and  $\nabla_j g(y) \leq \nabla_j g(x)$  otherwise.
2. Recall,  $\nabla_i g(x^{(*, C^{(t)})}) < 0$  only if  $i \in \mathcal{S}^* \setminus \mathcal{S}^{(t)}$ . So by 1., for  $x \in C^{(t)}$  s.t.  $x \leq x^{(*, C^{(t)})}$  it is  $\nabla_i g(x) < 0$  only if  $i \in \mathcal{S}^*$ :  
We can detect new coordinates!

## An inexact algorithm: Accelerated and Sparse PageRank (ASPR)

1. Because  $Q_{ij} \leq 0$  for  $i \neq j$ , for  $y = x - \Delta e_i$ , we have  $\forall j \neq i$ :  
 $\nabla_j g(y) \geq \nabla_j g(x)$  if  $\Delta > 0$  and  $\nabla_j g(y) \leq \nabla_j g(x)$  otherwise.
2. Recall,  $\nabla_i g(x^{(*, C^{(t)})}) < 0$  only if  $i \in \mathcal{S}^* \setminus \mathcal{S}^{(t)}$ . So by 1., for  $x \in C^{(t)}$  s.t.  $x \leq x^{(*, C^{(t)})}$  it is  $\nabla_i g(x) < 0$  only if  $i \in \mathcal{S}^*$ :  
We can detect new coordinates!
3. To ensure there exists such an  $i \in \mathcal{S}^* \setminus \mathcal{S}^{(t)}$ , get close to  $x^{(*, C^{(t)})}$  from below: optimize using accelerated projected gradient descent (APGD) to get close to  $x^{(*, C^{(t)})}$  and then move slightly towards  $\mathbf{0}$  to be  $\leq x^{(*, C^{(t)})}$ .

## An inexact algorithm: Accelerated and Sparse PageRank (ASPR)

1. Because  $Q_{ij} \leq 0$  for  $i \neq j$ , for  $y = x - \Delta e_i$ , we have  $\forall j \neq i$ :  
 $\nabla_j g(y) \geq \nabla_j g(x)$  if  $\Delta > 0$  and  $\nabla_j g(y) \leq \nabla_j g(x)$  otherwise.
2. Recall,  $\nabla_i g(x^{(*, C^{(t)})}) < 0$  only if  $i \in \mathcal{S}^* \setminus \mathcal{S}^{(t)}$ . So by 1., for  $x \in C^{(t)}$  s.t.  $x \leq x^{(*, C^{(t)})}$  it is  $\nabla_i g(x) < 0$  only if  $i \in \mathcal{S}^*$ :  
We can detect new coordinates!
3. To ensure there exists such an  $i \in \mathcal{S}^* \setminus \mathcal{S}^{(t)}$ , get close to  $x^{(*, C^{(t)})}$  from below: optimize using accelerated projected gradient descent (APGD) to get close to  $x^{(*, C^{(t)})}$  and then move slightly towards  $\mathbf{0}$  to be  $\leq x^{(*, C^{(t)})}$ .
4. **Lemma.** Optimizing with accuracy  $\hat{\varepsilon}_t = \varepsilon \cdot \frac{\alpha^2}{2(1+|\mathcal{S}^{(t)}|)L^2}$  to get  $\bar{x}^{(t+1)}$  and reducing  $x^{(t+1)} \leftarrow \max\{\mathbf{0}, \bar{x}^{(t+1)} - \delta_t \mathbf{1}\}$  for  $\delta_t = \sqrt{\frac{\varepsilon \alpha}{(1+|\mathcal{S}^{(t)}|)L^2}}$ , we either expand  $\mathcal{S}^{(t)}$  using 2. with  $x^{(t+1)}$ , or  $x^{(t+1)}$  is an  $\varepsilon$ -minimizer.

## An inexact algorithm: Accelerated and Sparse PageRank (ASPR)

1. Because  $Q_{ij} \leq 0$  for  $i \neq j$ , for  $y = x - \Delta e_i$ , we have  $\forall j \neq i$ :  
 $\nabla_j g(y) \geq \nabla_j g(x)$  if  $\Delta > 0$  and  $\nabla_j g(y) \leq \nabla_j g(x)$  otherwise.
2. Recall,  $\nabla_i g(x^{*, C^{(t)}}) < 0$  only if  $i \in \mathcal{S}^* \setminus \mathcal{S}^{(t)}$ . So by 1., for  $x \in C^{(t)}$  s.t.  $x \leq x^{*, C^{(t)}}$  it is  $\nabla_i g(x) < 0$  only if  $i \in \mathcal{S}^*$ :  
We can detect new coordinates!
3. To ensure there exists such an  $i \in \mathcal{S}^* \setminus \mathcal{S}^{(t)}$ , get close to  $x^{*, C^{(t)}}$  from below: optimize using accelerated projected gradient descent (APGD) to get close to  $x^{*, C^{(t)}}$  and then move slightly towards  $\mathbf{0}$  to be  $\leq x^{*, C^{(t)}}$ .
4. **Lemma.** Optimizing with accuracy  $\hat{\varepsilon}_t = \varepsilon \cdot \frac{\alpha^2}{2(1+|\mathcal{S}^{(t)}|)L^2}$  to get  $\bar{x}^{(t+1)}$  and reducing  $x^{(t+1)} \leftarrow \max\{\mathbf{0}, \bar{x}^{(t+1)} - \delta_t \mathbf{1}\}$  for  $\delta_t = \sqrt{\frac{\varepsilon \alpha}{(1+|\mathcal{S}^{(t)}|)L^2}}$ , we either expand  $\mathcal{S}^{(t)}$  using 2. with  $x^{(t+1)}$ , or  $x^{(t+1)}$  is an  $\varepsilon$ -minimizer.
5. The lemma above is proven by showing that if the global gap is  $> \varepsilon$ , then one step of gradient descent reduces the function value more than what it can be reduced in  $C^{(t)}$ .

# Accelerated and Sparse PageRank (ASPR) algorithm

## Theorem

The iterates of ASPR satisfy:

1.  $x_i^{(*,t)} > 0$  if and only if  $i \in S^{(t-1)}$ . Also,  $\nabla_i g(x^{(*,t)}) = 0$  if  $i \in S^{(t-1)}$ .
2. It is  $x^{(t)} \leq x^{(*,t)} \leq x^*$  and  $x^{(*,t-1)} \leq x^{(*,t)}$ .
3.  $S^{(t-1)} \subsetneq S^{(t)} \stackrel{\text{def}}{=} S^{(t-1)} \cup \{i \in [n] \mid \nabla_i g(x^{(t)}) < 0\} \subseteq S^*$ , or  $x^{(t)}$  is an  $\varepsilon$ -minimizer of  $g$ .

- ▶ APGD only needs gradients restricted to  $C^{(t)}$ , costing  $\mathcal{O}(\widetilde{\text{vol}}(S^*))$  each. Then, it uses a full gradient to find the new good coordinates, costing  $\mathcal{O}(\text{vol}(S^*))$ . It is done at most  $|S^*|$  times.
- ▶ All new good coordinates are incorporated to  $S^{(t)}$  unlike for CDPR.
- ▶ Time complexity  $\widetilde{\mathcal{O}}(|S^*| \widetilde{\text{vol}}(S^*) \sqrt{\frac{L}{\alpha}} + |S^*| \text{vol}(S^*))$  and space complexity  $\mathcal{O}(|S^*|)$ .



# Variants

- ▶ **Lemma.**  $S \subseteq [n]$ . If  $x$  is s.t.  $x_j = 0$  if  $j \notin S$  and  $\nabla_j g(x) \leq 0$  if  $j \in S$ , then for any  $i \notin S$  s.t.  $\nabla_i g(x) < 0$ , it is  $i \in \mathcal{S}^*$ .
- ▶ **Variant:** During APGD's execution, one can compute the full gradient from time to time to check the condition, expand  $S^{(t)}$ , and restart.

# Variants

- ▶ **Lemma.**  $S \subseteq [n]$ . If  $x$  is s.t.  $x_j = 0$  if  $j \notin S$  and  $\nabla_j g(x) \leq 0$  if  $j \in S$ , then for any  $i \notin S$  s.t.  $\nabla_i g(x) < 0$ , it is  $i \in \mathcal{S}^*$ .
- ▶ **Variant:** During **APGD**'s execution, one can compute the full gradient from time to time to check the condition, expand  $S^{(t)}$ , and restart.
- ▶ Interestingly, if we compute a full gradient at every iteration, **ASPR** is not better than **CDPR**, up to constants and logs, in the regime in which ISTA is not better, up to constants and logs.

# Variants

- ▶ **Lemma.**  $S \subseteq [n]$ . If  $x$  is s.t.  $x_j = 0$  if  $j \notin S$  and  $\nabla_j g(x) \leq 0$  if  $j \in S$ , then for any  $i \notin S$  s.t.  $\nabla_i g(x) < 0$ , it is  $i \in S^*$ .
- ▶ **Variant:** During **APGD**'s execution, one can compute the full gradient from time to time to check the condition, expand  $S^{(t)}$ , and restart.
- ▶ Interestingly, if we compute a full gradient at every iteration, **ASPR** is not better than **CDPR**, up to constants and logs, in the regime in which ISTA is not better, up to constants and logs.
- ▶ **Lemma.** If we observe  $\nabla_j g(x) \leq 0$  for all  $i \in S^{(t)}$ , it is  $x \leq x^{(*,t+1)} \leq x^*$ .
- ▶ **Variant:** With such an  $x$ , we can update the feasible set:  $C \leftarrow C \cap \{y \mid y \geq x\}$ .

# Variants

- ▶ **Lemma.**  $S \subseteq [n]$ . If  $x$  is s.t.  $x_j = 0$  if  $j \notin S$  and  $\nabla_j g(x) \leq 0$  if  $j \in S$ , then for any  $i \notin S$  s.t.  $\nabla_i g(x) < 0$ , it is  $i \in \mathcal{S}^*$ .
- ▶ **Variant:** During **APGD**'s execution, one can compute the full gradient from time to time to check the condition, expand  $S^{(t)}$ , and restart.
- ▶ Interestingly, if we compute a full gradient at every iteration, **ASPR** is not better than **CDPR**, up to constants and logs, in the regime in which ISTA is not better, up to constants and logs.
- ▶ **Lemma.** If we observe  $\nabla_j g(x) \leq 0$  for all  $i \in S^{(t)}$ , it is  $x \leq x^{(*,t+1)} \leq x^*$ .
- ▶ **Variant:** With such an  $x$ , we can update the feasible set:  $C \leftarrow C \cap \{y \mid y \geq x\}$ .
- ▶ **Variant:** Using the (unconstrained) Conjugate Gradients algorithm (CG) instead of **APGD**, the guarantee improves. And we can forgo the knowledge of the strong convexity constant  $\alpha$ .

# Comparisons

Method	Time complexity	Space complexity
ISTA [FRS+19]	$\tilde{\mathcal{O}}(\text{vol}(\mathcal{S}^*) \frac{L}{\alpha})$	$\mathcal{O}( \mathcal{S}^* )$
CDPR ( <b>Ours</b> )	$\mathcal{O}( \mathcal{S}^* ^3 +  \mathcal{S}^*  \text{vol}(\mathcal{S}^*))$	$\mathcal{O}( \mathcal{S}^* ^2)$
ASPR ( <b>Ours</b> )	$\tilde{\mathcal{O}}( \mathcal{S}^*  \widetilde{\text{vol}}(\mathcal{S}^*) \sqrt{\frac{L}{\alpha}} +  \mathcal{S}^*  \text{vol}(\mathcal{S}^*))$	$\mathcal{O}( \mathcal{S}^* )$
CASPR ( <b>Ours</b> )	$\tilde{\mathcal{O}}( \mathcal{S}^*  \widetilde{\text{vol}}(\mathcal{S}^*) \min \left\{ \sqrt{\frac{L}{\alpha}},  \mathcal{S}^*  \right\} +  \mathcal{S}^*  \text{vol}(\mathcal{S}^*))$	$\mathcal{O}( \mathcal{S}^* )$