

Semi-Proximal Mirror Prox for Nonsmooth Composite Minimization

Niao He

ISyE, Georgia Tech

Joint work with Zaid Harchaoui (NYU & Inria)

ISMP, Pittsburgh, July 17th, 2015

Outline

1 Background

2 Key Components

Semi-structured Variational Inequality

Composite Mirror Prox

Composite Conditional Gradient

3 The Semi-Proximal Mirror Prox Algorithm

4 Experiments

First-order methods for composite minimization

$$\min_{x \in X} f(x) + h(x)$$

f and h are convex, f is smooth, h is simple.

- **(Acc)Proximal gradient methods** (when h proximal-friendly)
 - Proximal operator:

$$\text{prox}_h(\eta) = \operatorname{argmin}_{x \in X} \left\{ \frac{1}{2} \|x - \eta\|_2^2 + h(x) \right\}$$

- For example, when $h(x) = \|x\|_1$, reduces to soft thresholding.
 - Worst complexity bound for first-order oracles is $\mathcal{O}(1/\sqrt{\epsilon})$.
- **Conditional gradient methods** (when h is LMO-friendly)
 - (Composite) linear minimization oracles(LMO):

$$\text{LMO}_h(\eta) = \operatorname{argmin}_{x \in X} \{ \langle \eta, x \rangle + h(x) \}$$

- For example, when $h(x) = \|x\|_{\text{nuc}}$ or $\delta_{\|x\|_{\text{nuc}} \leq 1}(x)$, reduces to computing top pair of singular vectors.
 - Worst (also optimal) complexity bound for LMOs is $\mathcal{O}(1/\epsilon)$.

Structured nonsmooth composite minimization

$$\min_{x \in X} f(x) + h(x)$$

f enjoys the *Fenchel-type representation*

$$f(x) = \max_{z \in Z} \{\langle x, Az \rangle - g(z)\}$$

for some convex function g and convex compact set Z .

Existing LMO-based algorithms are limited to

- **Nesterov's smoothing technique when (Z, g) are simple**
 - Run conditional gradient with the smooth approximation f^μ , where $|f^\mu - f| \leq O(\mu)$. [Lan'13] [PHM'14] [LZ'14]
- **Dual approach when Z admits favorable geometry**
 - Run first-order method on the dual problem [CJN'13] [JN'14]

Worst (also optimal) complexity bound for LMOs is $\mathcal{O}(1/\epsilon^2)$.

Problem of interest

$$\min_{x \in X} f(x) + h_1(x) + h_2(x)$$

- f (perhaps nonsmooth) admits **Fenchel-type representation**

$$f(x) = \max_{z \in Z} \Phi(x, z)$$

(**without** assuming Z to be simple or have favorable geometry)

- h_1 (perhaps nonsmooth) is **proximal-friendly**

$$\text{prox}_{h_1}(\eta) = \operatorname{argmin}_{x \in X} \left\{ \frac{1}{2} \|x - \eta\|_2^2 + h_1(x) \right\}$$

- h_2 (perhaps nonsmooth) is **LMO-friendly**

$$\text{LMO}_{h_2}(\eta) = \operatorname{argmin}_{x \in X} \left\{ \langle \eta, x \rangle + h_2(x) \right\}$$

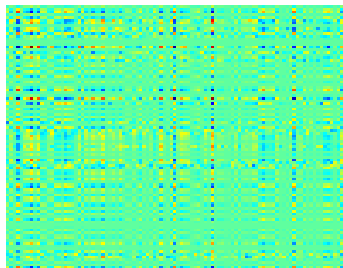
Example: social network link prediction

$$\min_{x \in \mathbf{R}^{n \times n}} \frac{1}{|E|} \sum_{(i,j) \in E} \underbrace{\max(1 - b_{ij}x_{ij}, 0)}_{f(x)} + \underbrace{\mu_1 \|x\|_1}_{h_1(x)} + \underbrace{\mu_2 \|x\|_{\text{nuc}}}_{h_2(x)}$$

- $f(x)$ = empirical hinge loss
- $h_1(x) = \ell_1$ norm \rightarrow sparsity
- $h_2(x) =$ nuclear norm \rightarrow low-rank



Social Network



Low-rank and Sparse Matrix

The goal and the cure

To design an algorithm that solves the above problem

$$\min_{x \in X} f(x) + h_1(x) + h_2(x)$$

by taking advantages and leveraging

- saddle point representation of f
- proximal operators of h_1
- linear minimization oracles of h_2

We propose **Semi-Proximal Mirror-Prox** that blends *proximal gradient* and *conditional gradient* and achieves best of both worlds

- optimal $\mathcal{O}(1/\epsilon)$ complexity bound in first-order oracles
- optimal $\mathcal{O}(1/\epsilon^2)$ complexity bound in linear minimization oracles

Course of action: saddle point reformulation

$$\min_{x \in X} f(x) + h_1(x) + h_2(x)$$

$$\Updownarrow$$

$$\min_{x, y \in X} \max_{z \in Z} \{\Phi(x, z) + h_1(x) + h_2(y) : x = y\}$$

$$\Updownarrow \text{ (with proper } \rho > 0 \text{)}$$

$$\min_{x, y \in X} \max_{z \in Z} \Phi(x, z) + h_1(x) + h_2(y) + \rho \|x - y\|_*$$

$$\Updownarrow$$

$$\min_{\substack{x, y \in X \\ \tau_1, \tau_2}} \max_{\substack{z \in Z \\ \|w\| \leq 1}} \Phi(x, z) + \rho \langle x - y, w \rangle + \tau_1 + \tau_2 \\ \tau_1 \geq h_1(x), \tau_2 \geq h_2(y)$$

Denote

$$\theta = [u; v], u = [x, y, z, w], v = [\tau_1, \tau_2]$$

$$\Theta = \{\theta = [u; v] : x, y \in X, z \in Z$$

$$\|w\| \leq 1, \tau_1 \geq h_1(x), \tau_2 \geq h_2(y)\}$$

Variational Inequality VI(Θ, F)

$$\langle F(\theta), \theta - \theta_* \rangle \geq 0, \forall \theta \in \Theta$$

where

$$F(\theta) = [F_u(u); F_v]$$

$$F_u(u) = \begin{bmatrix} \nabla_x \Phi(x, z) + \rho w \\ -\rho w \\ -\nabla_z \Phi(x, z) \\ \rho(y - x) \end{bmatrix}$$

$$F_v = [1; 1]$$

Semi-structured variational inequality

$$VI(\Theta, F) : \quad \langle F(\theta), \theta - \theta_* \rangle \geq 0, \forall \theta \in \Theta$$

- Decomposition of $F(\theta = [u, v]) = [F_u(u); F_v]$:

$$\|F_u(u) - F_u(u')\| \leq L\|u - u'\| \text{ and } F_v \text{ is constant}$$

- Decomposition of $\Theta = \Theta_1 \times \Theta_2$:

$$\theta = [\theta_1 = [u_1, v_1]; \theta_2 = [u_2, v_2]]$$

- “Computation-friendly” proximal mapping on Θ_1 :

$$\text{Prox}(\eta_1, \alpha) := \min_{\theta_1 \in \Theta_1} \{ \omega_1(u_1) + \langle \eta_1, u_1 \rangle + \alpha \langle F_{v_1}, v_1 \rangle \}.$$

for some distance generating function $\omega_1(\cdot)$.

- “Computation-friendly” linear minimization on Θ_2 :

$$\text{LMO}(\eta_2, \alpha) := \min_{\theta_2 \in \Theta_2} \{ \langle \eta_2, u_2 \rangle + \alpha \langle F_{v_2}, v_2 \rangle \}.$$

Composite Mirror Prox with inexact prox-mappings

$$VI(\Theta, F) : \quad \langle F(\theta), \theta - \theta_* \rangle \geq 0, \forall \theta \in \Theta$$

$$F(\theta) = [F_u(u); F_v]$$

- ϵ -prox-mapping on entire Θ : let $\omega(u)$ be the d.s.g. on u -space, for any $\xi = [\eta, \zeta]$ and $\theta_0 = [u_0, v_0] \in \Theta$,

$$P_{\theta_0}^\epsilon(\xi) = \{[\hat{u}, \hat{v}] \in \Theta : \langle \eta + \omega'(\hat{u}) - \omega'(u_0), \hat{u} - u \rangle + \langle \zeta, \hat{v} - v \rangle \leq \epsilon, \forall [u; v] \in \Theta\}$$

Inexact Composite Mirror Prox

$$(a) \quad \theta^1 := [u^1; v^1] \in \Theta;$$

$$(b) \quad \hat{\theta}^t := [\hat{u}^t; \hat{v}^t] \in P_{\theta^t}^{\epsilon_t}(\gamma_t [F_u(u^t); F_v])$$

$$(c) \quad \theta^{t+1} := [u^{t+1}; v^{t+1}] \in P_{\theta^t}^{\epsilon_t}(\gamma_t [F_u(\hat{u}^t); F_v])$$

- When $\epsilon_t = 0$, reduces to the composite Mirror Prox [HJN'13] with exact prox-mapping

$$P_{\theta^t}(\gamma_t [F_u(u^t); F_v]) = \underset{\theta=[u;v] \in \Theta}{\text{Argmin}} \{ D_\omega(u, u^t) + \langle \gamma_t F_u(u^t), u \rangle + \langle \gamma_t F_v, v \rangle \}.$$

Composite Mirror Prox with inexact prox-mappings

Theorem [HH'15]

Let $\bar{\theta}_T = \left[\sum_{t=1}^T \gamma_t \right]^{-1} \sum_{t=1}^T \gamma_t \hat{\theta}^t$ be generated by the inexact composite Mirror Prox algorithm by setting $0 < \gamma_t \leq \frac{1}{L}$ and inexactness $\epsilon_t \geq 0$. Let Ω be the bound for prox function. Then

$$\epsilon_{\text{VI}}(\bar{\theta}_T | \Theta, F) := \sup_{\theta \in \Theta} \langle F(\theta), \bar{\theta}_T - \theta \rangle \leq \frac{\Omega + 2 \sum_{t=1}^T \epsilon_t}{\sum_{t=1}^T \gamma_t}.$$

Immediate observations

- When $\epsilon_t \equiv 0$, one has at least $\epsilon_{\text{VI}}(\bar{\theta}_T | \Theta, F) \leq \frac{\Omega L}{T}$.
- When $\epsilon_t \equiv \frac{\Omega}{2T}$, one has at least $\epsilon_{\text{VI}}(\bar{\theta}_T | \Theta, F) \leq \frac{2\Omega L}{T}$.
- When $\epsilon_t = \frac{\Omega}{2t}$, one has at least $\epsilon_{\text{VI}}(\bar{\theta}_T | \Theta, F) \leq \frac{\Omega L(1 + \ln T)}{T}$.

Inexact prox-mapping via LMO routines

- Recall the subproblem of prox-mapping in composite Mirror Prox

$$P_{\theta^t}(\gamma_t[F_u(u^t); F_v]) = \underset{\theta=[u;v] \in \Theta}{\operatorname{Argmin}} \{D_\omega(u, u^t) + \langle \gamma_t F_u(u^t), u \rangle + \langle \gamma_t F_v, v \rangle\}.$$

reduces to solve the **smooth semi-linear problem**:

$$\min_{\theta=[u;v] \in \Theta} \{\phi^+(u, v) = \phi(u) + \langle c, v \rangle\}.$$

- Assume we have at disposal only the linear minimization oracles

$$\operatorname{LMO}(\eta) := \min_{\theta=[u;v] \in \Theta} \{\langle \eta, u \rangle + \langle c, v \rangle\}.$$

Composite Conditional Gradient

$$\min_{\theta=[u,v] \in \Theta} \{ \phi^+(u, v) = \phi(u) + \langle c, v \rangle \}.$$

Composite Conditional Gradient $CCG(\Theta, \phi, c; \epsilon)$

$$(a) \quad \theta^1 := [u^1; v^1] \in \Theta;$$

$$(b) \quad \hat{\theta}^t := [\hat{u}^t; \hat{v}^t] = LMO(\nabla \phi(u^t))$$

compute $\delta_t = \langle \nabla \phi^+(\theta^t), \theta^t - \hat{\theta}^t \rangle$, return if $\delta_t \leq \epsilon$

$$(c) \quad \theta^{t+1} := [u^{t+1}; v^{t+1}] \text{ s.t. } \phi^+(\theta^{t+1}) \leq \phi^+(\theta^t + \frac{2}{t+1}(\hat{\theta}^t - \theta^t))$$

- Choices of θ^{t+1} in (c):
 - Simplest: $\theta^{t+1} = \theta^t + \gamma_t(\hat{\theta}^t - \theta^t)$ with $\gamma_t = 2/(t+1)$.
 - Line search: $\theta^{t+1} = \operatorname{argmin}_{0 \leq \gamma \leq 1} \phi^+(\gamma \hat{\theta}^t + (1-\gamma)\theta^t)$.
 - Fully-corrective search: $\theta_{t+1} = \operatorname{argmin}_{\theta \in \operatorname{conv}(\theta^1, \dots, \theta^t)} \phi^+(\theta)$
- Use the dual gap δ_t to check the condition for ϵ -prox mapping.

Composite Conditional Gradient

Theorem [HJN'13, HH'15]

Assume convex function ϕ is (κ, L_0) -smooth for some $1 < \kappa \leq 2$,

$$\phi(u') \leq \phi(u) + \langle \nabla \phi(u), u' - u \rangle + \frac{L_0}{\kappa} \|u' - u\|^\kappa \quad \forall u, u';$$

Let D the $\|\cdot\|$ -diameter of $P_u \Theta$. Then

$$\phi^+(x^t) - \min_{x \in X} \phi^+(x) \leq \frac{2L_0 D^\kappa}{\kappa(3-\kappa)} \left(\frac{2}{t+1} \right)^{\kappa-1}.$$

In addition, the accuracy certificates δ_t satisfy

$$\min_{1 \leq s \leq t} \delta_s \leq O(1) L_0 D^\kappa \left(\frac{2}{t+1} \right)^{\kappa-1}.$$

- When $\phi(u) = \frac{1}{2} \|u - u_0\|_2^2$, the worst complexity bound for LMO calls is $\mathcal{O}(1/\epsilon)$.

Returning to semi-structured variational inequality

$$VI(\Theta, F) : \quad \langle F(\theta), \theta - \theta_* \rangle \geq 0, \forall \theta \in \Theta$$

- Decomposition of $F(\theta = [u, v]) = [F_u(u); F_v]$:

$$\|F_u(u) - F_u(u')\| \leq L\|u - u'\| \text{ and } F_v \text{ is constant}$$

- Decomposition of $\Theta = \Theta_1 \times \Theta_2$:

$$\theta = [\theta_1 = [u_1, v_1]; \theta_2 = [u_2, v_2]]$$

- “Computation-friendly” proximal mapping on Θ_1 :

$$\text{Prox}(\eta_1, \alpha) := \min_{\theta_1 \in \Theta_1} \{ \omega_1(u_1) + \langle \eta_1, u_1 \rangle + \alpha \langle F_{v_1}, v_1 \rangle \}.$$

for some distance generating function $\omega_1(\cdot)$.

- “Computation-friendly” linear minimization on Θ_2 :

$$\text{LMO}(\eta_2, \alpha) := \min_{\theta_2 \in \Theta_2} \{ \langle \eta_2, u_2 \rangle + \alpha \langle F_{v_2}, v_2 \rangle \}.$$

Working horse: Semi-Proximal Mirror-Prox

Semi-Proximal Mirror-Prox

Input: stepsizes $\gamma_t > 0$, accuracies $\epsilon_t \geq 0$, $t = 1, 2, \dots$

[1] Initialize $\theta^1 = [\theta_1^1; \theta_2^1] \in \Theta$, where $\theta_1^1 = [u_1^1; v_1^1]; \theta_2^1 = [u_2^1; v_2^1]$.

for $t = 1, 2, \dots, T$ **do**

[2] Compute $\hat{\theta}^t = [\hat{\theta}_1^t; \hat{\theta}_2^t]$ that

$$\hat{\theta}_1^t := [\hat{u}_1^t; \hat{v}_1^t] = \text{Prox}(\gamma_t F_{u_1}(u_1^t) - \omega_1'(u_1^t), \gamma_t)$$

$$\hat{\theta}_2^t := [\hat{u}_2^t; \hat{v}_2^t] = \text{CCG}(X_2, \omega_2(\cdot) + \langle \gamma_t F_{u_2}(u_2^t) - \omega_2'(u_2^t), \cdot \rangle, \gamma_t F_{v_2}; \epsilon_t)$$

[3] Compute $\theta^{t+1} = [\theta_1^{t+1}; \theta_2^{t+1}]$ that

$$\theta_1^{t+1} := [u_1^{t+1}; v_1^{t+1}] = \text{Prox}(\gamma_t F_{u_1}(\hat{u}_1^t) - \omega_1'(u_1^t), \gamma_t)$$

$$\theta_2^{t+1} := [u_2^{t+1}; v_2^{t+1}] = \text{CCG}(X_2, \omega_2(\cdot) + \langle \gamma_t F_{u_2}(\hat{u}_2^t) - \omega_2'(u_2^t), \cdot \rangle, \gamma_t F_{v_2}; \epsilon_t)$$

end for

Output: $\bar{\theta}_T := [\bar{u}_T; \bar{v}_T] = (\sum_{t=1}^T \gamma_t)^{-1} \sum_{t=1}^T \gamma_t \hat{\theta}^t$

Interpretating "semi-proximal"

- When Θ_2 is singleton, we get *full-proximal setup*, the algorithm reduces to composite Mirror Prox.
- When Θ_1 is singleton, we get *full-LMO setup*, the algorithm serves as a viable alternative to existing LMO-based algorithms via smoothing technique or dual approach for solving nonsmooth composite minimization.
- In the gray zone in between, we get *semi-proximal setup*, the algorithm is the first of its kind to solve semi-structured variational inequality, which covers the aforementioned problem of interest.

Complexity

Theorem

For the outlined algorithm to return an ϵ -solution to the $VI(\Theta, F)$, the total number of Mirror Prox steps required does not exceed

$$T = O\left(\frac{L\Omega}{\epsilon}\right)$$

and the total number LMO calls required does not exceed

$$\mathcal{N} = O(1) \left(\frac{L_0 L^\kappa D^\kappa}{\epsilon^\kappa}\right)^{\frac{1}{\kappa-1}} \Omega.$$

In particular, if we use $\omega_2(\cdot) = \frac{1}{2}\|x_2\|^2$, then $\mathcal{N} = O(1)L^2 D^2 \Omega / \epsilon^2$.

Immediate observation

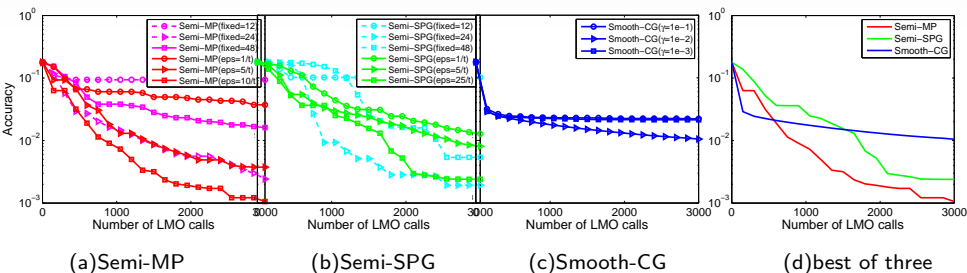
- optimal complexity bound $O(1/\epsilon)$ in first-order oracles
- optimal complexity bound $O(1/\epsilon^2)$ in linear minimization oracles

Numerical experiments

- Several types of problems
 - Matrix completion
 - Robust collaborative filtering
 - Link prediction
- Several alternative algorithms
 - **Semi-MP**: Semi-Proximal Mirror-Prox
 - **Smooth-CG**: Smoothing conditional gradient [PHM'14]
 - **Semi-SPG**: Smoothing proximal gradient with inexact prox-mappings via LMO routines
 - **Semi-LPADMM**: Linearized preconditioned ADMM with inexact prox-mappings via LMO routines
- Several choices of inexactness
 - Use decaying error $\epsilon_t = \frac{c}{t}$
 - Use fixed inner CG steps

Matrix completion: ℓ_2 -fit + nuclear norm

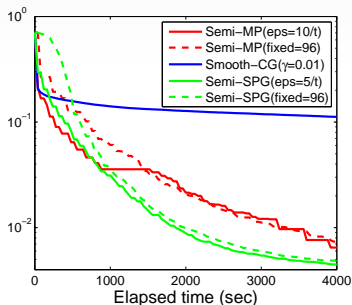
$$\min_{x \in \mathbf{R}^{m \times n}} \|P_{\Omega}x - b\|_2 + \lambda \|x\|_{\text{nuc}}$$



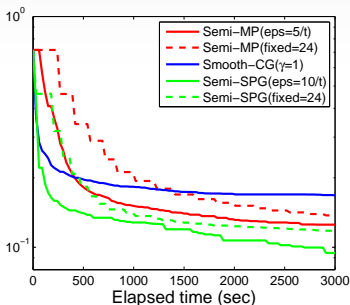
Synthetic dataset(1024×1024): optimality gap vs the LMO calls.

Robust collaborative filtering: ℓ_1 -loss + nuclear norm

$$\min_{x \in \mathbf{R}^{m \times n}} \frac{1}{|E|} \sum_{(i,j) \in E} |x_{ij} - b_{ij}| + \lambda \|x\|_{\text{nuc}}$$



(a) MovieLens 100K(943 × 1682)

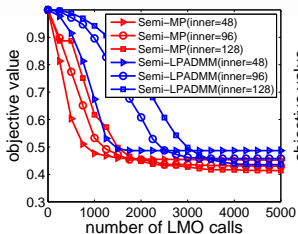
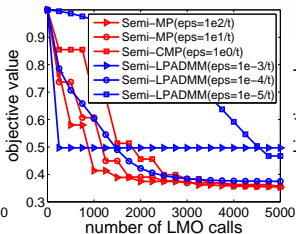
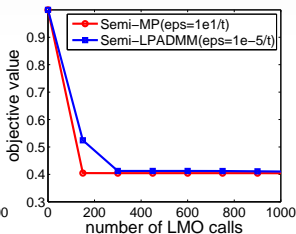


(b) MovieLens 1M(3952 × 6040)

MovieLens dataset: objective function vs elapsed time.

Link prediction: hinge loss + ℓ_1 -norm + nuclear norm

$$\min_{x \in \mathbf{R}^{m \times n}} \frac{1}{|E|} \sum_{(i,j) \in E} \max(1 - b_{ij}x_{ij}, 0) + \lambda_1 \|x\|_1 + \lambda_2 \|x\|_{\text{nuc}}$$

(a) Wikivote(1024 nodes)
fixed inner steps(b) Wikivote(1024 nodes)
decaying $\epsilon_t = c/t$ (c) Wikivote(7118 nodes)
decaying $\epsilon_t = c/t$

Link prediction on Wikivote: objective function value against the LMO calls.

Conclusion and perspectives

- We present a new algorithm, Semi-Proximal Mirror-Prox , that solves a broad class of composite minimization and leverages
 - saddle point representation of f
 - proximal operators of h_1
 - linear minimization oracles of h_2
- Theoretically, the algorithm achieves
 - optimal complexity bound $\mathcal{O}(1/\epsilon)$ in first-order oracles
 - optimal complexity bound $\mathcal{O}(1/\epsilon^2)$ in linear minimization oracles
- Practically, the algorithm shows advantages in comparison to competing methods in various large-scale applications.
- The algorithm is readily extensible to online/stochastic settings.

Thank you!

References

- [HH'15] *Semi-proximal Mirror-Prox for nonsmooth composite minimization*, N. He and Z. Harchaoui, 2015.
- [HJN'15] *Mirror Prox algorithm for multi-term composite minimization and semi-separable problems*, N. He, A. Juditsky, and A. Nemirovski, 2015.
- [HJN'13] *Conditional gradient algorithms for norm-regularized smooth convex optimization*, Z. Harchaoui, A. Juditsky, and A. Nemirovski, 2013.
- [Lan'13] *The complexity of large-scale convex programming under a linear optimization oracle*, G. Lan, 2013.
- [LZ'14] *Conditional gradient sliding for convex optimization*, G. Lan and Y. Zhou, 2014.
- [PHM'14] *A smoothing approach for composite conditional gradient with nonsmooth loss*, F. Pierucci, Z. Harchaoui, and J. Malick, 2014.