

Recall the setting:

$$v^* = \inf_{\theta \in \mathbb{R}^d} \{ \varphi(\theta) \triangleq F(\theta) + r(\theta) \} \quad - (P)$$

where $F: \mathbb{R}^d \rightarrow \mathbb{R}$ is p -weakly convex and $r = \mathbb{1}_C$ for some closed convex set C .

We assume that F has a finite-sum structure:

$$F(\theta) = \sum_{i=1}^n F_i(\theta).$$

Projected stochastic subgradient method (PSSM)

a) Sample $\{1, \dots, n\}$ uniformly at random to get ξ_t

$$\Pr[\xi_t = i] = \frac{1}{n}; \quad i = 1, \dots, n.$$

b) $\theta^{t+1} \leftarrow \Pi_C(\theta^t - \alpha_t s(\theta^t, \xi_t))$, $s(\theta^t, \xi_t) \in \partial F_{\xi_t}(\theta^t)$

To analyze PSSM, we need a measure that can monitor its progress.

As discussed, the norm of the Moreau envelope of φ is a good

candidate.

Definition: Given $\varphi: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ and $\lambda > 0$, define

$$\text{prox}_{\lambda\varphi}(\theta) = \underset{\gamma \in \mathbb{R}^d}{\text{argmin}} \left\{ \varphi(\gamma) + \frac{1}{2\lambda} \|\gamma - \theta\|_2^2 \right\} \quad (\text{proximal map})$$

$$\varphi_\lambda(\theta) = \min_{\gamma \in \mathbb{R}^d} \left\{ \varphi(\gamma) + \frac{1}{2\lambda} \|\gamma - \theta\|_2^2 \right\} \quad (\text{Moreau envelope})$$

We say that $\bar{\theta} \in \mathbb{R}^d$ is an (ε, δ) -approximate stationary point of

(P) if $\bar{\theta} \in \{ \gamma : \text{dist}(\theta, \partial\varphi(\gamma)) \leq \varepsilon \} + \delta B(0, 1)$.

With the above setup, let us analyze PSSM under the following

assumptions:

$$A) \mathbb{E}_\xi[s(\theta, \xi)] \in \partial F(\theta), \quad \forall \theta \in U \text{ open}, \quad U \supseteq \text{dom}(r)$$

A) $\mathbb{E}_{\xi} [s(\theta, \xi)] \in \partial F(\theta)$, $\forall \theta \in U$ open, $U \supseteq \text{dom}(r)$
(unbiasedness)

B) $\mathbb{E}_{\xi} [\|s(\theta, \xi)\|_2^2] \leq L^2$, $\forall \theta \in \text{dom}(r)$ (variance control)

C) After T iterations, output θ^{t^*} , where $t^* \in \{0, 1, \dots, T\}$ is sampled according to

$$\Pr [t^* = t] = \frac{\alpha_t}{\sum_{j=0}^T \alpha_j}.$$

Theorem: (Iteration Complexity of PSSM)

For any $\bar{\rho} > \rho$,

$$\mathbb{E} \left[\underbrace{\varphi_{\frac{1}{\bar{\rho}}}}_{\lambda}(\theta^{t+1}) \right] \leq \mathbb{E} \left[\varphi_{\frac{1}{\bar{\rho}}}(\theta^t) \right] - \frac{\alpha_t(\bar{\rho} - \rho)}{\bar{\rho}} \mathbb{E} \left[\|\nabla \varphi_{\frac{1}{\bar{\rho}}}(\theta^t)\|_2^2 \right] + \frac{\alpha_t^2 \bar{\rho} L^2}{2},$$

("sufficient decrease" in terms of $\varphi_{\frac{1}{\bar{\rho}}}$)

so that

$$\mathbb{E} \left[\|\nabla \varphi_{\frac{1}{\bar{\rho}}}(\theta^{t^*})\|_2^2 \right] \leq \frac{\bar{\rho}}{\bar{\rho} - \rho} \frac{(\varphi_{\frac{1}{\bar{\rho}}}(\theta^0) - v^*) + \frac{\bar{\rho} L^2}{2} \sum_{t=0}^T \alpha_t^2}{\sum_{t=0}^T \alpha_t}.$$

In particular, if we take $\bar{\rho} = 2\rho$, $\alpha_t = \sqrt{\frac{\Delta}{\rho L^2 (T+1)}}$ for some $\Delta > \varphi_{\frac{1}{2\rho}}(\theta^0) - v^*$, then

$$\mathbb{E} \left[\|\nabla \varphi_{\frac{1}{2\rho}}(\theta^{t^*})\|_2^2 \right] \leq 4 \sqrt{\frac{\rho \Delta L^2}{T+1}}.$$

Proof: Let $\hat{\theta}^t = \text{prox}_{\frac{1}{\bar{\rho}} \varphi}(\theta^t)$

$$\mathbb{E}_t \left[\varphi_{\frac{1}{\bar{\rho}}}(\theta^{t+1}) \right] \leq \mathbb{E}_t \left[\varphi(\hat{\theta}^t) + \frac{\bar{\rho}}{2} \|\hat{\theta}^t - \theta^{t+1}\|_2^2 \right]$$

↑
Expectation with info up to time t

$$= \mathbb{E}_t \left[F(\hat{\theta}^t) + \frac{\bar{\rho}}{2} \|\hat{\theta}^t - \theta^{t+1}\|_2^2 \right]$$

(Recall: $\varphi(\theta) = F(\theta) + \mathbb{1}_C(\theta)$. Note that $\hat{\theta}^t \in C$, so $\mathbb{1}_C(\hat{\theta}^t) = 0$)

$$= F(\hat{\theta}^t) + \frac{\bar{\rho}}{2} \mathbb{E}_t \left[\underbrace{\|\pi_C(\theta^t - \alpha_t s(\theta^t, \xi_t)) - \pi_C(\hat{\theta}^t)\|_2^2}_{\hat{\theta}^t \dots \hat{\theta}^t \dots} \right]$$

$$= F(\theta^t) + \frac{\rho}{2} \mathbb{E}_t \left[\left\| \underbrace{\Pi_C(\theta^t - \alpha_t s(\theta^t, \xi_t))}_{\theta^{t+1}} - \underbrace{\Pi_C(\hat{\theta}^t)}_{\hat{\theta}^t} \right\|_2^2 \right] \quad (\because \hat{\theta}^t \in C)$$

$$\leq F(\hat{\theta}^t) + \frac{\bar{\rho}}{2} \mathbb{E}_t \left[\left\| \theta^t - \hat{\theta}^t - \alpha_t s(\theta^t, \xi_t) \right\|_2^2 \right]$$

(by the non-expansiveness of Π_C . for any $\theta, \gamma \in \mathbb{R}^d$,

$$\left\| \Pi_C(\gamma) - \Pi_C(\theta) \right\|_2 \leq \left\| \gamma - \theta \right\|_2$$

$$= \underbrace{F(\hat{\theta}^t) + \frac{\bar{\rho}}{2} \left\| \theta^t - \hat{\theta}^t \right\|_2^2}_{\text{}} + \bar{\rho} \alpha_t \mathbb{E}_t \left[(\hat{\theta}^t - \theta^t)^T s(\theta^t, \xi_t) \right] + \frac{\bar{\rho} \alpha_t^2}{2} \mathbb{E}_t \left[\left\| s(\theta^t, \xi_t) \right\|_2^2 \right]$$

$$= \underbrace{\varphi_{\frac{1}{\bar{\rho}}}(\theta^t)}_{\hookrightarrow \hat{\theta}^t = \text{prox}_{\frac{1}{\bar{\rho}}} \varphi(\theta^t)} + \bar{\rho} \alpha_t (\hat{\theta}^t - \theta^t)^T \underbrace{\mathbb{E}_t \left[s(\theta^t, \xi_t) \right]}_{\substack{\triangleq s^t \in \partial F(\theta^t) \\ \text{by assumption (A)}}} + \frac{\bar{\rho} \alpha_t^2}{2} \underbrace{\mathbb{E}_t \left[\left\| s(\theta^t, \xi_t) \right\|_2^2 \right]}_{\leq L^2 \text{ by assumption (B)}}$$

$$\leq \varphi_{\frac{1}{\bar{\rho}}}(\theta^t) + \bar{\rho} \alpha_t \left[F(\hat{\theta}^t) - F(\theta^t) + \frac{\rho}{2} \left\| \theta^t - \hat{\theta}^t \right\|_2^2 \right] + \frac{\bar{\rho} \alpha_t^2 L^2}{2}$$

(by weak convexity of F)

Since $\theta \mapsto F(\theta) + \frac{\bar{\rho}}{2} \left\| \theta - \theta^t \right\|_2^2$ is strongly convex with parameter $\bar{\rho} - \rho$, we have

$$F(\theta^t) - F(\hat{\theta}^t) - \frac{\rho}{2} \left\| \theta^t - \hat{\theta}^t \right\|_2^2 = \underbrace{\left(F(\theta^t) + \frac{\bar{\rho}}{2} \left\| \theta^t - \theta^t \right\|_2^2 \right)}_{=0} - \underbrace{\left(F(\hat{\theta}^t) + \frac{\bar{\rho}}{2} \left\| \theta^t - \hat{\theta}^t \right\|_2^2 \right)}_{\substack{+ \frac{\bar{\rho} - \rho}{2} \left\| \theta^t - \hat{\theta}^t \right\|_2^2 \\ - \underbrace{\theta^{t \top} (\theta^t - \hat{\theta}^t)}}_{\text{note that}}}$$

$$\geq (\bar{\rho} - \rho) \left\| \theta^t - \hat{\theta}^t \right\|_2^2 = \frac{\bar{\rho} - \rho}{\bar{\rho}^2} \left\| \nabla \varphi_{\frac{1}{\bar{\rho}}}(\theta^t) \right\|_2^2$$

↑
use strong convexity of $\theta \mapsto F(\theta) + \frac{\bar{\rho}}{2} \left\| \theta - \theta^t \right\|_2^2$

note that
 $\theta \in \partial \left(F(\gamma) + \frac{\bar{\rho}}{2} \left\| \gamma - \theta^t \right\|_2^2 \right) \Big|_{\gamma = \hat{\theta}^t}$
 Since by definition,
 $\hat{\theta}^t = \text{prox}_{\frac{1}{\bar{\rho}}} \varphi(\theta^t)$.

Putting the pieces together yield the first inequality in the theorem