

The basic PGM for solving

$$\hat{\theta} \in \underset{\theta \in \mathbb{R}^d}{\operatorname{Argmin}} \left\{ F(\theta) \triangleq \underbrace{\ell(\theta)}_{\substack{\text{Convex} \\ \nabla \ell \text{ L-Lipschitz}}} + \underbrace{R(\theta)}_{\text{Convex}} \right\} \quad (*)$$

is given by

$$\begin{aligned} \theta^{k+1} &= \operatorname{prox}_{\alpha_k R} \left(\theta^k - \alpha_k \nabla \ell(\theta^k) \right) \\ &= \underset{\gamma \in \mathbb{R}^d}{\operatorname{Argmin}} \left\{ \frac{1}{2} \left\| \theta^k - \alpha_k \nabla \ell(\theta^k) - \gamma \right\|_2^2 + \alpha_k R(\gamma) \right\}. \end{aligned}$$

Proposition 1: Suppose that $\{\alpha_k\}$ satisfy $0 < \underline{\alpha} < \alpha_k < \bar{\alpha} < \frac{1}{L}$ for some $\underline{\alpha}, \bar{\alpha}$. Assuming $(*)$ is solvable,

(a) (sufficient decrease) There exists a constant $K_1 > 0$ s.t.

$$F(\theta^k) - F(\theta^{k+1}) \geq K_1 \|\theta^k - \theta^{k+1}\|_2^2$$

(b) (Safeguard; relative error) There exists a constant $K_2 > 0$ s.t.

$$\|E(\theta^k)\|_2 \leq K_2 \|\theta^k - \theta^{k+1}\|_2, \quad E(\theta) = \operatorname{prox}_R(\theta - \nabla \ell(\theta)) - \theta$$

Consequences

$$(1) \quad F(\theta^k) \rightarrow \hat{v} \triangleq F(\hat{\theta}),$$

(2) Every accumulation point of $\{\theta^k\}$ is optimal for $(*)$.

Goal: Determine the convergence rate of PGM.

Need a measure of the progress of the method.

Proposition 2: Let Θ^* be the optimal solution set of $(*)$, assumed to be non-empty (thus, Convex and closed). Under the setting of Proposition 1,

(c) (Cost-to-go estimate) There exists a $K_3 > 0$ s.t.

$$F(\theta^{k+1}) - \hat{v} \leq \kappa_3 \left[\text{dist}(\theta^k, \Theta)^2 + \|\theta^{k+1} - \theta^k\|_2^2 \right],$$

where $\text{dist}(\theta, \Theta) = \inf_{\gamma \in \Theta} \|\theta - \gamma\|_2$.

Proof: Let $\bar{\theta}^k = \Pi_{\Theta}(\theta^k)$. By definition, $F(\theta^{k+1}) = \mathcal{L}(\theta^{k+1}) + R(\theta^{k+1})$

and since $\theta^{k+1} = \text{prox}_{\alpha_k R}(\theta^k - \alpha_k \nabla \mathcal{L}(\theta^k))$,

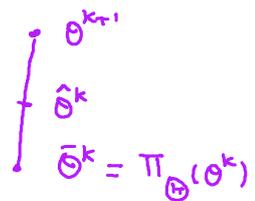
$$\begin{aligned} & \frac{1}{2} \|\theta^{k+1} - \theta^k + \alpha_k \nabla \mathcal{L}(\theta^k)\|_2^2 + \alpha_k R(\theta^{k+1}) \\ & \leq \frac{1}{2} \|\bar{\theta}^k - \theta^k + \alpha_k \nabla \mathcal{L}(\theta^k)\|_2^2 + \alpha_k R(\bar{\theta}^k) \end{aligned}$$

$$\Rightarrow R(\theta^{k+1}) - R(\bar{\theta}^k) + \nabla \mathcal{L}(\theta^k)^T (\theta^{k+1} - \bar{\theta}^k) \quad \leftarrow$$

$$\leq \frac{1}{2\alpha_k} \|\bar{\theta}^k - \theta^k\|_2^2 \leq \frac{1}{2\underline{\alpha}} \text{dist}(\theta^k, \Theta)^2.$$

On the other hand, by the mean-value theorem, there exists $\hat{\theta}^k \in [\bar{\theta}^k, \theta^{k+1}]$ s.t.

$$\mathcal{L}(\theta^{k+1}) - \mathcal{L}(\bar{\theta}^k) = \nabla \mathcal{L}(\hat{\theta}^k)^T (\theta^{k+1} - \bar{\theta}^k)$$



It follows that

$$\begin{aligned} F(\theta^{k+1}) - \hat{v} &= F(\theta^{k+1}) - F(\bar{\theta}^k) \\ &= \mathcal{L}(\theta^{k+1}) - \mathcal{L}(\bar{\theta}^k) + R(\theta^{k+1}) - R(\bar{\theta}^k) \\ &= \nabla \mathcal{L}(\theta^k)^T (\theta^{k+1} - \bar{\theta}^k) + R(\theta^{k+1}) - R(\bar{\theta}^k) \quad \leftarrow \\ &\quad + (\nabla \mathcal{L}(\hat{\theta}^k) - \nabla \mathcal{L}(\theta^k))^T (\theta^{k+1} - \bar{\theta}^k) \\ &\leq \frac{1}{2\underline{\alpha}} \text{dist}(\theta^k, \Theta)^2 + L \cdot \|\hat{\theta}^k - \theta^k\|_2 \cdot \underbrace{\|\theta^{k+1} - \bar{\theta}^k\|_2}_{\text{green underline}} \end{aligned}$$

Now, note that

$$\|\hat{\theta}^k - \theta^k\|_2 \leq \|\theta^{k+1} - \theta^k\|_2 + \|\theta^{k+1} - \hat{\theta}^k\|_2 \quad (\text{triangle inequality})$$

$$\leq \| \theta^{k+1} - \theta^k \|_2 + \| \theta^{k+1} - \bar{\theta}^k \|_2,$$

$$\| \theta^{k+1} - \bar{\theta}^k \|_2 \leq \| \theta^{k+1} - \theta^k \|_2 + \| \theta^k - \bar{\theta}^k \|_2$$

$$= \| \theta^{k+1} - \theta^k \|_2 + \text{dist}(\theta^k, \Theta)$$

Hence,

$$\| \hat{\theta}^k - \theta^k \|_2 \| \theta^{k+1} - \bar{\theta}^k \|_2 \leq 2 \left[\text{dist}(\theta^k, \Theta) + \| \theta^{k+1} - \theta^k \|_2 \right]^2$$

$$\leq 4 \left[\text{dist}(\theta^k, \Theta)^2 + \| \theta^{k+1} - \theta^k \|_2^2 \right]$$

Since $\forall a, b \in \mathbb{R}, (a+b)^2 \leq 2(a^2 + b^2)$

Observe that $\text{dist}(\theta^k, \Theta)$ is not easy to estimate. To circumvent this difficulty, let us assume the following:

(A4) (Local Error Bound Condition)

For any $v \geq \hat{v}$, there exist $\mu, \epsilon > 0$ s.t.

$$\text{dist}(\theta, \Theta) \leq \mu \cdot \|E(\theta)\|_2 \quad (E(\theta) = \text{prox}_{\mathbb{R}}(\theta - \nabla f(\theta)) - \theta)$$

for any $\theta \in \mathbb{R}^d$ satisfying $F(\theta) \leq v, \|E(\theta)\|_2 \leq \epsilon$.

Remarks

(1) (A4) says when θ is "close" to Θ (as determined by the conditions $F(\theta) \leq v, \|E(\theta)\|_2 \leq \epsilon$), the measure $\|E(\theta)\|_2$ is a good surrogate of $\text{dist}(\theta, \Theta)$. However, if μ is independent of v and $\epsilon = +\infty$, then we say that the error bound is global.

(2) (A4) is a property of the problem and is algorithm independent.

Using Propositions 1 and 2 and (A4), we can get the convergence rate of PGM. Indeed, for Θ^k close to

Θ in the sense of (A4), we have

$$\begin{aligned}
 F(\Theta^{k+1}) - \hat{v} &\stackrel{(A4)}{\leq} K_4 \left[\|\epsilon(\Theta^k)\|_2^2 + \|\Theta^k - \Theta^{k+1}\|_2^2 \right] \\
 &\stackrel{(b)}{\leq} K_5 \|\Theta^k - \Theta^{k+1}\|_2^2 \\
 &\stackrel{(a)}{\leq} K_6 \left[(F(\Theta^k) - \hat{v}) - (F(\Theta^{k+1}) - \hat{v}) \right]
 \end{aligned}$$

Hence,

$$F(\Theta^{k+1}) - \hat{v} \leq \underbrace{\frac{K_6}{1+K_6}}_{< 1} (F(\Theta^k) - \hat{v})$$

$\Rightarrow \{F(\Theta^k) - \hat{v}\}$ converges to 0 geometrically

Similarly, by (a),

$$\begin{aligned}
 \|\Theta^k - \Theta^{k+1}\|_2^2 &\leq \frac{1}{K_1} \left[(F(\Theta^k) - \hat{v}) - \underbrace{(F(\Theta^{k+1}) - \hat{v})}_{\geq 0} \right] \\
 &\leq \frac{1}{K_1} (F(\Theta^k) - \hat{v})
 \end{aligned}$$

Thus,

$$\begin{aligned}
 \text{dist}(\Theta^k, \Theta) &\stackrel{(A4)}{\leq} \mu \cdot \|\epsilon(\Theta^k)\|_2 \stackrel{(b)}{\leq} \mu' \cdot \|\Theta^k - \Theta^{k+1}\|_2 \\
 &\leq \mu'' \cdot \sqrt{F(\Theta^k) - \hat{v}}
 \end{aligned}$$

This shows $\{\Theta^k\}$ approaches Θ geometrically.