Notes for Lecture 7

Scribe: Tingkai Jia

1 Projected Subgradient Descent with Polyak's Stepsize

Lemma 1. Projected subgradient update obeys

$$\|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2^2 \le \|\mathbf{x}_t - \mathbf{x}_*\|_2^2 + \eta_t^2 \|\mathbf{g}_t\|_2^2 - 2\eta_t (f(\mathbf{x}_t) - f(\mathbf{x}_*)).$$

Proof. It follows that

$$\begin{aligned} \|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2^2 &= \|\mathcal{P}_{\mathcal{C}}(\mathbf{x}_t - \eta_t \mathbf{g}_t) - \mathbf{x}_*\|_2^2 \\ &\leq \|\mathbf{x}_t - \eta_t \mathbf{g}_t - \mathbf{x}_*\|_2^2 \\ &= \|\mathbf{x}_t - \mathbf{x}_*\|_2^2 + \eta_t^2 \|\mathbf{g}_t\|_2^2 - 2\eta_t \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}_* \rangle \\ &\leq \|\mathbf{x}_t - \mathbf{x}_*\|_2^2 + \eta_t^2 \|\mathbf{g}_t\|_2^2 - 2\eta_t (f(\mathbf{x}_t) - f(\mathbf{x}_*)), \end{aligned}$$

where the last inequality uses the convexity

$$f(\mathbf{x}_*) \geq f(\mathbf{x}_t) + \langle \mathbf{g}_t, \mathbf{x}_* - \mathbf{x}_t \rangle.$$

Definition 1 (polyak's stepsize). By treating the RHS of Inequality in Lemma 1 as a quadratic function with respect to η_t , we obtain a step size by minimizing this function

$$\eta_t = \frac{f(\mathbf{x}_t) - f(\mathbf{x}_*)}{\|\mathbf{g}_t\|_2^2}.$$

1.1 Example: Projection onto Intersection of Convex Sets

Example 1. Let C_1, C_2 be closed convex sets and suppose $C_1 \cap C_2 \neq \emptyset$,

$$minimize_x \quad max\{dist_{\mathcal{C}_1}(\mathbf{x}), dist_{\mathcal{C}_2}(\mathbf{x})\}$$

where $\operatorname{dist}_{\mathcal{C}}(\mathbf{x}) := \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{x} - \mathbf{z}\|_{2}$.

For this problem, the subgradient method of polyak's stepsize will act as

$$\mathbf{x}_{t+1} = \mathcal{P}_{\mathcal{C}_1}(\mathbf{x}_t), \quad \mathbf{x}_{t+2} = \mathcal{P}_{\mathcal{C}_2}(\mathbf{x}_{t+1}).$$

Proof. First we consider the subgradient, it follows that

$$\mathbf{g}_t \in \partial \mathrm{dist}_{\mathcal{C}_i}(\mathbf{x}_t)$$

where $i = \arg \max_{i=1,2} \operatorname{dist}_{\mathcal{C}_i}(\mathbf{x}_t)$. If $\operatorname{dist}_{\mathcal{C}_i}(\mathbf{x}_t) \neq 0$, we have

$$\mathbf{g}_t = \nabla dist_{\mathcal{C}_i}(\mathbf{x}_t) = \frac{\mathbf{x}_t - \mathcal{P}_{\mathcal{C}_i}(\mathbf{x}_t)}{\|\mathbf{x}_t - \mathcal{P}_{\mathcal{C}_i}(\mathbf{x}_t)\|_2}.$$

Then the polyak's stepsize is shown as

$$\eta_t = \frac{\operatorname{dist}_{\mathcal{C}_i}(\mathbf{x}_t) - 0}{\|\mathbf{g}_t\|_2^2} = \|\mathbf{x}_t - \mathcal{P}_{\mathcal{C}_i}(\mathbf{x}_t)\|_2.$$

Adopting polyak's stepsize, we arrive at

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \mathbf{g}_t = \mathbf{x}_t - \|\mathbf{x}_t - \mathcal{P}_{\mathcal{C}_i}(\mathbf{x}_t)\|_2 \frac{\mathbf{x}_t - \mathcal{P}_{\mathcal{C}_i}(\mathbf{x}_t)}{\|\mathbf{x}_t - \mathcal{P}_{\mathcal{C}_i}(\mathbf{x}_t)\|_2} = \mathcal{P}_{\mathcal{C}_i}(\mathbf{x}_t).$$

1.2 Convergence Rate with Polyak's Stepsize

Theorem 1. Suppose f is convex and L-Lipschitz continuous. Then the projected subgradient method with Polyak's stepsize obeys

$$f_{\text{best,t}} - f(\mathbf{x}_*) \le \frac{L \|\mathbf{x}_0 - \mathbf{x}_*\|_2}{\sqrt{t+1}}$$

Proof. With Lemma 1 and substituting η_t , we obtain

$$(f(\mathbf{x}_t) - f(\mathbf{x}_*))^2 \le [\|\mathbf{x}_t - \mathbf{x}_*\|_2^2 - \|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2^2] \|\mathbf{g}_t\|_2^2$$

$$\le [\|\mathbf{x}_t - \mathbf{x}_*\|_2^2 - \|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2^2] L^2.$$

Applying it resursively, we get

$$\sum_{k=0}^{t} (f(\mathbf{x}_{k}) - f(\mathbf{x}_{*}))^{2} \leq \left[\|\mathbf{x}_{0} - \mathbf{x}_{*}\|_{2}^{2} - \|\mathbf{x}_{t+1} - \mathbf{x}_{*}\|_{2}^{2} \right] L^{2}$$

$$(t+1)(f_{\text{best,t}} - f(\mathbf{x}_{*}))^{2} \leq \left[\|\mathbf{x}_{0} - \mathbf{x}_{*}\|_{2}^{2} - \|\mathbf{x}_{t+1} - \mathbf{x}_{*}\|_{2}^{2} \right] L^{2}$$

$$(f_{\text{best,t}} - f(\mathbf{x}_{*}))^{2} \leq \frac{L^{2} \|\mathbf{x}_{0} - \mathbf{x}_{*}\|_{2}^{2}}{t+1}$$

which completes the proof.

2 Projected Subgradient Descent with Other Stepsizes

Lemma 2. Suppose f is convex and L-Lipschitz continuous. Then the projected subgradient update obeys

$$f_{\text{best,t}} - f(\mathbf{x}_*) \le \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 + L^2 \sum_{i=0}^t \eta_i^2}{2 \sum_{i=0}^t \eta_i}.$$

Proof. Using Lemma 1 and summing it recursively, we obtain

$$2\sum_{i=0}^{t} \eta_{i}(f(\mathbf{x}_{i}) - f(\mathbf{x}_{*})) \leq \|\mathbf{x}_{0} - \mathbf{x}_{*}\|_{2}^{2} - \|\mathbf{x}_{t+1} - \mathbf{x}_{*}\|_{2}^{2} + \sum_{i=0}^{t} \eta_{i}^{2} \|\mathbf{g}_{i}\|_{2}^{2}$$

$$2(f_{\text{best,t}} - f(\mathbf{x}_{*})) \sum_{i=0}^{t} \eta_{i} \leq \|\mathbf{x}_{0} - \mathbf{x}_{*}\|_{2}^{2} - \|\mathbf{x}_{t+1} - \mathbf{x}_{*}\|_{2}^{2} + \sum_{i=0}^{t} \eta_{i}^{2} \|\mathbf{g}_{i}\|_{2}^{2}$$

$$2(f_{\text{best,t}} - f(\mathbf{x}_{*})) \sum_{i=0}^{t} \eta_{i} \leq \|\mathbf{x}_{0} - \mathbf{x}_{*}\|_{2}^{2} - \|\mathbf{x}_{t+1} - \mathbf{x}_{*}\|_{2}^{2} + \sum_{i=0}^{t} \eta_{i}^{2} L^{2}$$

$$f_{\text{best,t}} - f(\mathbf{x}_{*}) \leq \frac{\|\mathbf{x}_{0} - \mathbf{x}_{*}\|_{2}^{2} + L^{2} \sum_{i=0}^{t} \eta_{i}^{2}}{2\sum_{i=0}^{t} \eta_{i}},$$

thus we complete the proof.

2.1 Convergence with $1/\sqrt{t+1}$ Stepsize

Considering the inequality in Lemma 2, we aim to make its RHS approach zero as the subgradient method update, which means $\sum_{i=0}^{t} \eta_i^2 < \infty$ and $\sum_{i=0}^{t} \eta_i \to \infty$. Now we can consider $\eta_t = \frac{1}{\sqrt{t+1}}$.

Theorem 2. Suppose f is convex and L-Lipschitz continuous. Then the projected subgradient method with $\eta_t = \frac{1}{\sqrt{t+1}}$ obeys

$$f_{\text{best,t}} - f(\mathbf{x}_*) \lesssim \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 + L^2}{\sqrt{t}}.$$

Proof. With the fact that $\frac{2}{\sqrt{t+1}+\sqrt{t+2}} \le \frac{1}{\sqrt{t+1}} \le \frac{2}{\sqrt{t+\sqrt{t+1}}}$, we can get a lower and upper bound in $\sum_{i=0}^t \eta_i$,

$$\sum_{k=0}^{t} \frac{2}{\sqrt{k+1} + \sqrt{k+2}} \le \sum_{i=0}^{t} \eta_i \le \sum_{k=0}^{t} \frac{2}{\sqrt{k} + \sqrt{k+1}}$$
$$2\sum_{k=0}^{t} (\sqrt{k+2} - \sqrt{k+1}) \le \sum_{i=0}^{t} \eta_i \le 2\sum_{k=0}^{t} (\sqrt{k+1} - \sqrt{k})$$
$$2(\sqrt{t+2} - 1) \le \sum_{i=0}^{t} \eta_i \le 2(\sqrt{t+1}),$$

and consider sequence $\frac{1}{k}$, the upper bound of its sum is $\log t + 1$, now we get

$$f_{\text{best,t}} - f(\mathbf{x}_*) \le \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 + L^2 \sum_{i=0}^t \eta_i^2}{2 \sum_{i=0}^t \eta_i}$$

$$\le \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 + L^2 (\log t + 1)}{4(\sqrt{t+2} - 1)}$$

$$\lesssim \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 + L^2 \log t}{\sqrt{t}}.$$

Through this approach, we find that the convergence contains a $\log t$ in the numerator. now, we attempt to eliminate this term.

Note that $\sum_{k=\lceil \frac{t}{2} \rceil}^t \frac{1}{k} \approx \log t - \log \lceil \frac{t}{2} \rceil \le \log 3$ and $\sum_{k=\lceil \frac{t}{2} \rceil}^t \frac{1}{\sqrt{k}} \approx 2\sqrt{t} - 2\sqrt{\lceil \frac{t}{2} \rceil} = (2 - \sqrt{2})\sqrt{t}$. Thus, we modify the inequality in Lemma 2 and obtain

$$2 \sum_{i=\lceil \frac{t}{2} \rceil}^{t} \eta_{i}(f(\mathbf{x}_{i}) - f(\mathbf{x}_{*})) \leq \|\mathbf{x}_{\lceil \frac{t}{2} \rceil} - \mathbf{x}_{*}\|_{2}^{2} - \|\mathbf{x}_{t+1} - \mathbf{x}_{*}\|_{2}^{2} + L^{2} \sum_{i=\lceil \frac{t}{2} \rceil}^{t} \eta_{i}^{2}$$

$$2(f_{\text{best},t} - f(\mathbf{x}_{*}) \sum_{i=\lceil \frac{t}{2} \rceil}^{t} \eta_{i} \leq \|\mathbf{x}_{\lceil \frac{t}{2} \rceil} - \mathbf{x}_{*}\|_{2}^{2} - \|\mathbf{x}_{t+1} - \mathbf{x}_{*}\|_{2}^{2} + L^{2} \sum_{i=\lceil \frac{t}{2} \rceil}^{t} \eta_{i}^{2}$$

$$2(f_{\text{best},t} - f(\mathbf{x}_{*}) \sum_{i=\lceil \frac{t}{2} \rceil}^{t} \eta_{i} \leq \|\mathbf{x}_{0} - \mathbf{x}_{*}\|_{2}^{2} + L^{2} \sum_{i=\lceil \frac{t}{2} \rceil}^{t} \eta_{i}^{2}$$

$$f_{\text{best},t} - f(\mathbf{x}_{*}) \leq \frac{\|\mathbf{x}_{0} - \mathbf{x}_{*}\|_{2}^{2} + L^{2} \log 3}{2(2 - \sqrt{2})\sqrt{t}}$$

$$f_{\text{best},t} - f(\mathbf{x}_{*}) \lesssim \frac{\|\mathbf{x}_{0} - \mathbf{x}_{*}\|_{2}^{2} + L^{2} \log 3}{\sqrt{t}},$$

which we finish the proof.

3 Strongly Convex and Lipschitz Problems

Theorem 3. Let f be -strongly convex and L-Lipschitz continuous over C. If $\eta_t \equiv \eta = \frac{2}{\mu(t+1)}$, then

$$f_{\text{best,t}} - f(\mathbf{x}_*) \le \frac{2L^2}{\mu(t+1)}$$

Proof. Consider strongly convex situation in Lemma 1, we have

$$\begin{aligned} \|\mathbf{x}_{t+1} - \mathbf{x}_{*}\|_{2}^{2} &= \|\mathcal{P}_{\mathcal{C}}(\mathbf{x}_{t} - \eta_{t}\mathbf{g}_{t}) - \mathbf{x}_{*}\|_{2}^{2} \\ &\leq \|\mathbf{x}_{t} - \eta_{t}\mathbf{g}_{t} - \mathbf{x}_{*}\|_{2}^{2} \\ &= \|\mathbf{x}_{t} - \mathbf{x}_{*}\|_{2}^{2} + \eta_{t}^{2}\|\mathbf{g}_{t}\|_{2}^{2} - 2\eta_{t}\langle\mathbf{g}_{t}, \mathbf{x}_{t} - \mathbf{x}_{*}\rangle \\ &\leq (1 - \mu\eta_{t})\|\mathbf{x}_{t} - \mathbf{x}_{*}\|_{2}^{2} + \eta_{t}^{2}\|\mathbf{g}_{t}\|_{2}^{2} - 2\eta_{t}(f(\mathbf{x}_{t}) - f(\mathbf{x}_{*})), \end{aligned}$$

where the last inequality uses the μ -strongly convexity

$$f(\mathbf{x}_*) \ge f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x}_* - \mathbf{x}_t \rangle + \frac{\mu}{2} ||\mathbf{x}_* - \mathbf{x}_t||_2^2.$$

Since $\eta_t \equiv \eta = \frac{2}{\mu(t+1)}$, we have

$$f(\mathbf{x}_{t}) - f(\mathbf{x}_{*}) \leq \frac{\mu(t-1)}{4} \|\mathbf{x}_{t} - \mathbf{x}_{*}\|_{2}^{2} - \frac{\mu(t+1)}{4} \|\mathbf{x}_{t+1} - \mathbf{x}_{*}\|_{2}^{2} + \frac{1}{\mu(t+1)} \|\mathbf{g}_{t}\|_{2}^{2}$$

$$t(f(\mathbf{x}_{t}) - f(\mathbf{x}_{*})) \leq \frac{\mu(t-1)}{4} \|\mathbf{x}_{t} - \mathbf{x}_{*}\|_{2}^{2} - \frac{\mu(t+1)}{4} \|\mathbf{x}_{t+1} - \mathbf{x}_{*}\|_{2}^{2} + \frac{t}{\mu(t+1)} \|\mathbf{g}_{t}\|_{2}^{2}$$

$$\leq \frac{\mu(t-1)}{4} \|\mathbf{x}_{t} - \mathbf{x}_{*}\|_{2}^{2} - \frac{\mu(t+1)}{4} \|\mathbf{x}_{t+1} - \mathbf{x}_{*}\|_{2}^{2} + \frac{1}{\mu} \|\mathbf{g}_{t}\|_{2}^{2}.$$

Summing over all iterations before t, we get

$$\sum_{k=0}^{t} k(f(\mathbf{x}_k) - f(\mathbf{x}_*)) \le 0 - \frac{\mu t(t+1)}{4} \|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2^2 + \frac{1}{\mu} \sum_{i=0}^{t} \|\mathbf{g}_k\|_2^2$$

$$\le \frac{tL^2}{\mu},$$

which means

$$f_{\text{best,t}} - f(\mathbf{x}_*) \le \frac{tL^2}{\mu \sum_{k=0}^t k} \le \frac{2L^2}{\mu(t+1)}.$$

Thus we finish the proof.