

Convergence analysis

Convex and Lipschitz problems

$$\begin{array}{ll} \text{minimize}_x & f(x) \\ \text{subject to} & x \in \mathcal{C} \\ & \text{s.t.} \end{array}$$

- f is convex and Lipschitz continuous *with respect to*
 - φ is ρ -strongly convex w.r.t. a certain norm $\|\cdot\|$
 - $\|g\|_* \leq L_f$ for any subgradient $g \in \partial f(x)$ at any point x , where $\|\cdot\|_*$ is the dual norm of $\|\cdot\|$

$$\|g\|_* = \sup \{ \langle g, x \rangle \mid \|x\| \leq 1 \}$$

Convergence analysis

Theorem 5.3

Suppose f is convex and Lipschitz continuous (in the sense that $\|g\|_* \leq L_f$ for any subgradient g of f) on \mathcal{C} . Suppose φ is ρ -strongly convex w.r.t. $\|\cdot\|$. Then

$$f^{\text{best},t} - f^{\text{opt}} \leq \frac{\sup_{\mathbf{x} \in \mathcal{C}} D_\varphi(\mathbf{x}, \mathbf{x}^0) + \frac{L_f^2}{2\rho} \sum_{k=0}^t \eta_k^2}{\sum_{k=0}^t \eta_k}$$

- If $\eta_t = \frac{\sqrt{2\rho R}}{L_f} \frac{1}{\sqrt{t}}$ with $R := \sup_{\mathbf{x} \in \mathcal{C}} D_\varphi(\mathbf{x}, \mathbf{x}^0)$, then

$$f^{\text{best},t} - f^{\text{opt}} \leq O\left(\frac{L_f \sqrt{R} \log t}{\sqrt{\rho} \sqrt{t}}\right)$$

- one can further remove the $\log t$ factor

Example: optimization over probability simplex

Suppose $\mathcal{C} = \Delta$ is the probability simplex, and pick $\mathbf{x}^0 = n^{-1}\mathbf{1}$

(1) set $\varphi(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|_2^2$, which is 1-strongly convex w.r.t. $\|\cdot\|_2$. Then

$$\sup_{\mathbf{x} \in \Delta} D_{\varphi}(\mathbf{x}, \mathbf{x}^0) = \sup_{\mathbf{x} \in \Delta} \frac{1}{2} \|\mathbf{x} - n^{-1}\mathbf{1}\|_2^2 = \sup_{\mathbf{x} \in \Delta} \frac{1}{2} \left(\|\mathbf{x}\|_2^2 - \frac{1}{n} \right) \leq \frac{1}{2}$$

Then Theorem 5.3 says $\|\mathbf{x}\|_2^2 - 2\langle \mathbf{x}, \frac{\mathbf{1}}{n} \rangle + \|\frac{\mathbf{1}}{n}\|_2^2$

$$f^{\text{best},t} - f^{\text{opt}} \leq O\left(L_{f,2} \frac{\log t}{\sqrt{t}}\right) = \|\mathbf{x}\|_2^2 - \frac{2}{n} + \frac{1}{n}$$

if any subgradient \mathbf{g} obeys $\|\mathbf{g}\|_2 \leq L_{f,2}$

$$= \|\mathbf{x}\|_2^2 - \frac{1}{n}$$

Example: optimization over probability simplex

Suppose $\mathcal{C} = \Delta$ is the probability simplex, and pick $\mathbf{x}^0 = n^{-1}\mathbf{1}$

(2) set $\phi(\mathbf{x}) = \textcolor{red}{+} \sum_{i=1}^n x_i \log x_i$, which is 1-strongly convex w.r.t. $\|\cdot\|_1$. Then

$$\begin{aligned} \sup_{\mathbf{x} \in \Delta} D_{\phi}(\mathbf{x}, \mathbf{x}^0) &= \sup_{\mathbf{x} \in \Delta} \text{KL}(\mathbf{x} \parallel \mathbf{x}^0) = \sup_{\mathbf{x} \in \Delta} \sum_{i=1}^n x_i \log x_i - \sum_{i=1}^n x_i \log \frac{1}{n} \\ &= \log n + \sup_{\mathbf{x} \in \Delta} \sum_{i=1}^n x_i \log x_i \leq \log n \end{aligned}$$

Then Theorem 5.3 says

$$f^{\text{best},t} - f^{\text{opt}} \leq O\left(L_{f,\infty} \sqrt{\log n} \frac{\log t}{\sqrt{t}}\right)$$

if any subgradient \mathbf{g} obeys $\|\mathbf{g}\|_{\infty} \leq L_{f,\infty}$

Example: optimization over probability simplex

Comparing these two choices and ignoring log terms, we have

$$\text{Euclidean: } O\left(\frac{L_{f,2}}{\sqrt{t}}\right) \quad \text{vs.} \quad \text{KL: } O\left(\frac{L_{f,\infty}}{\sqrt{t}}\right)$$

Since $\|\mathbf{g}\|_\infty \leq \|\mathbf{g}\|_2 \leq \sqrt{n}\|\mathbf{g}\|_\infty$, one has

$$\frac{1}{\sqrt{n}} \leq \frac{L_{f,\infty}}{L_{f,2}} \leq 1$$

and hence the KL version often yields much better performance

$$Loss_{MSE} = \min \frac{1}{2} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$N(0, 1)$$

$$p(y_i | x_i) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{(y_i - \hat{y}_i)^2}{2} \right\}$$

$$i = 1 \dots N$$

$$\max \log \left(\prod_{i=1}^N \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{(y_i - \hat{y}_i)^2}{2} \right\} \right)$$

$$\max -\frac{N}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$\min \frac{1}{2} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$\text{Binary: } \begin{cases} p(y_i=1|x_i) = p_i \quad \checkmark \\ p(y_i=0|x_i) = 1-p_i \quad \checkmark \end{cases}$$

$$p(y_i|x_i) = p_i^{y_i} (1-p_i)^{1-y_i} \quad \checkmark$$

$$\frac{N}{\log} \prod_{i=1}^N p_i^{y_i} (1-p_i)^{1-y_i}$$

$$\Downarrow$$

$$\min - \sum_{i=1}^N (y_i \log p_i + (1-y_i) \log (1-p_i))$$

$$\min - \sum_{i=1}^N \sum_{k=1}^K y_i^k \log p_i^k$$

Numerical example: robust regression

taken from Stanford EE364B

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) = \sum_{i=1}^m |\mathbf{a}_i^\top \mathbf{x} - b_i|$$

$$\text{subject to} \quad \mathbf{x} \in \Delta = \{\mathbf{x} \in \mathbb{R}_+^n \mid \mathbf{1}^\top \mathbf{x} = 1\}$$

with $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n \times n})$ and $b_i = \frac{a_{i,1} + a_{i,2}}{2} + \mathcal{N}(0, 10^{-2})$, $m = 20$,
 $n = 3000$

subgradient of f : $g \in \partial f$

$$g = \sum_{i=1}^m \text{sgn}(\mathbf{a}_i^\top \mathbf{x} - b_i) \mathbf{a}_i$$

$$\textcircled{1} \quad \varphi(x) = \frac{1}{2} \|x\|_2^2, \quad D\varphi(x, x_t) = \frac{1}{2} \|x - x^t\|_2^2$$

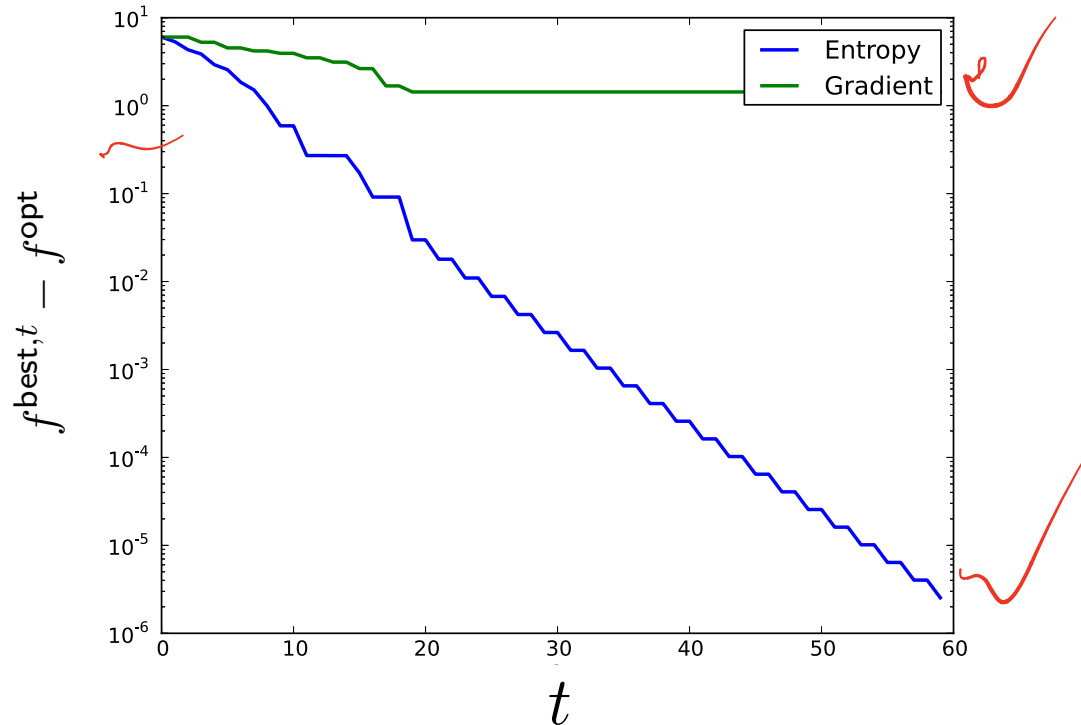
$$x^{t+1} = x^t - \eta_t g^t$$

$$\textcircled{2} \quad \varphi(x) = \sum_i x_i \log x_i, \quad D\varphi(x, x_t) = \sum_i x_i \log \frac{x_i}{x_i^t}$$

$$(x_i^{t+1}) = \frac{x_i^t \exp(-\eta_t g_i^t)}{\sum_j x_j^t \exp(-\eta_t g_j^t)} = \sum_i x_i \log x_i - \sum_i x_i \log x_i^t$$

Numerical example: robust regression

taken from Stanford EE364B



Fundamental inequality for mirror descent

Lemma 5.4

$$\eta_t \left(f(\mathbf{x}^t) - \underbrace{f^{\text{opt}}}_{\text{red circle}} \right) \leq D_\varphi(\mathbf{x}^*, \mathbf{x}^t) - D_\varphi(\mathbf{x}^*, \mathbf{x}^{t+1}) + \frac{\eta_t^2 L_f^2}{2\rho}$$

- $D_\varphi(\mathbf{x}^*, \mathbf{x}^t) - D_\varphi(\mathbf{x}^*, \mathbf{x}^{t+1})$ motivates us to form a telescopic sum later

Proof of Lemma 5.4

$$f(\mathbf{x}^t) - f(\mathbf{x}^*) \leq \langle \mathbf{g}^t, \mathbf{x}^t - \mathbf{x}^* \rangle \quad (\text{property of subgradient})$$

$$= \frac{1}{\eta_t} \langle \nabla \varphi(\mathbf{x}^t) - \nabla \varphi(\mathbf{y}^{t+1}), \mathbf{x}^t - \mathbf{x}^* \rangle \quad (\text{MD update rule})$$

$$= \frac{1}{\eta_t} \{ D_\varphi(\mathbf{x}^*, \mathbf{x}^t) + D_\varphi(\mathbf{x}^t, \mathbf{y}^{t+1}) - D_\varphi(\mathbf{x}^*, \mathbf{y}^{t+1}) \} \quad (\text{three point lemma})$$

$$\leq \frac{1}{\eta_t} \{ D_\varphi(\mathbf{x}^*, \mathbf{x}^t) + D_\varphi(\mathbf{x}^t, \mathbf{y}^{t+1}) - D_\varphi(\mathbf{x}^*, \mathbf{x}^{t+1}) - D_\varphi(\mathbf{x}^{t+1}, \mathbf{y}^{t+1}) \}$$

(Pythagorean)

$$= \frac{1}{\eta_t} \{ D_\varphi(\mathbf{x}^*, \mathbf{x}^t) - D_\varphi(\mathbf{x}^*, \mathbf{x}^{t+1}) \} + \frac{1}{\eta_t} \{ D_\varphi(\mathbf{x}^t, \mathbf{y}^{t+1}) - D_\varphi(\mathbf{x}^{t+1}, \mathbf{y}^{t+1}) \}$$

so we need to first bound the 2nd term of the last line

Proof of Lemma 5.4 (cont.)

We claim that

$$D_{\varphi}(\mathbf{x}^t, \mathbf{y}^{t+1}) - D_{\varphi}(\mathbf{x}^{t+1}, \mathbf{y}^{t+1}) \leq \frac{(\eta_t L_f)^2}{2\rho} \quad (5.6)$$

This gives

$$\eta_t (f(\mathbf{x}^t) - f(\mathbf{x}^*)) \leq \{D_{\varphi}(\mathbf{x}^*, \mathbf{x}^t) - D_{\varphi}(\mathbf{x}^*, \mathbf{x}^{t+1})\} + \frac{(\eta_t L_f)^2}{2\rho}$$

as claimed

Proof of Lemma 5.4 (cont.)

Finally, we justify (5.6):

$$\begin{aligned} & D_\varphi(\mathbf{x}^t, \mathbf{y}^{t+1}) - D_\varphi(\mathbf{x}^{t+1}, \mathbf{y}^{t+1}) \\ &= \varphi(\mathbf{x}^t) - \varphi(\mathbf{x}^{t+1}) - \langle \nabla \varphi(\mathbf{y}^{t+1}), \mathbf{x}^t - \mathbf{x}^{t+1} \rangle \quad \checkmark \\ &\leq \langle \nabla \varphi(\mathbf{x}^t), \mathbf{x}^t - \mathbf{x}^{t+1} \rangle - \frac{\rho}{2} \|\mathbf{x}^t - \mathbf{x}^{t+1}\|^2 - \langle \nabla \varphi(\mathbf{y}^{t+1}), \mathbf{x}^t - \mathbf{x}^{t+1} \rangle \\ &\hspace{25em} \text{(strong convexity of } \varphi) \\ &= \langle \nabla \varphi(\mathbf{x}^t) - \nabla \varphi(\mathbf{y}^{t+1}), \mathbf{x}^t - \mathbf{x}^{t+1} \rangle - \frac{\rho}{2} \|\mathbf{x}^t - \mathbf{x}^{t+1}\|_2^2 \\ &= \eta_t \langle \mathbf{g}^t, \mathbf{x}^t - \mathbf{x}^{t+1} \rangle - \frac{\rho}{2} \|\mathbf{x}^t - \mathbf{x}^{t+1}\|^2 \hspace{10em} \text{(MD update rule)} \\ &\leq \eta_t L_f \|\mathbf{x}^t - \mathbf{x}^{t+1}\| - \frac{\rho}{2} \|\mathbf{x}^t - \mathbf{x}^{t+1}\|^2 \hspace{10em} \text{(Cauchy-Schwarz)} \\ &\leq \frac{(\eta_t L_f)^2}{2\rho} \hspace{10em} \text{(optimize quadratic function in } \|\mathbf{x}^t - \mathbf{x}^{t+1}\|) \end{aligned}$$

$$\|\mathbf{x}^t - \mathbf{x}^{t+1}\|^2 = \frac{\eta_t L_f}{\rho}$$

Proof of Theorem 5.3

From Lemma 5.4, one has

$$\eta_k \left(f(\mathbf{x}^k) - f^{\text{opt}} \right) \leq D_\varphi(\mathbf{x}^*, \mathbf{x}^k) - D_\varphi(\mathbf{x}^*, \mathbf{x}^{k+1}) + \frac{\eta_k^2 L_f^2}{2\rho}$$

Taking this inequality for $k = 0, \dots, t$ and summing them up give

$$\begin{aligned} \sum_{k=0}^t \eta_k \left(f(\mathbf{x}^k) - f^{\text{opt}} \right) &\leq D_\varphi(\mathbf{x}^*, \mathbf{x}^0) - D_\varphi(\mathbf{x}^*, \mathbf{x}^{t+1}) + \frac{L_f^2 \sum_{k=0}^t \eta_k^2}{2\rho} \\ &\leq \sup_{\mathbf{x} \in \mathcal{C}} D_\varphi(\mathbf{x}, \mathbf{x}^0) + \frac{L_f^2 \sum_{k=0}^t \eta_k^2}{2\rho} \end{aligned}$$

This together with $f^{\text{best},t} - f^{\text{opt}} \leq \frac{\sum_{k=0}^t \eta_k (f(\mathbf{x}^k) - f^{\text{opt}})}{\sum_{k=0}^t \eta_k}$ concludes the proof

Reference

- [1] "*Problem complexity and method efficiency in optimization*," A. Nemirovski, D. Yudin, Wiley, 1983.
- [2] "*Mirror descent and nonlinear projected subgradient methods for convex optimization*," A. Beck, M. Teboulle, Operations Research Letters, 31(3), 2003.
- [3] "*Convex optimization: algorithms and complexity*," S. Bubeck, Foundations and trends in machine learning, 2015.
- [4] "*First-order methods in optimization*," A. Beck, Vol. 25, SIAM, 2017.
- [5] "*Mathematical optimization, MATH301 lecture notes*," E. Candes, Stanford.
- [6] "*Convex optimization, EE364B lecture notes*," S. Boyd, Stanford.

Reference

- [7] "*Matrix nearness problems with Bregman divergences*," I. Dhillon, J. Tropp, SIAM Journal on Matrix Analysis and Applications, 29(4), 2007.
- [8] "*Nonlinear Programming (2nd Edition)*," D. Bertsekas, Athena Scientific, 2016.