

Mirror Descent Part 2

Yanjie Ze. July 2021

Alternative forms of mirror descent

$$\text{original form: } x^{t+1} = \underset{x \in C}{\operatorname{argmin}} \left\{ \langle g^t, x - x^t \rangle + \frac{1}{\eta_t} D_\varphi(x, x^t) \right\}$$

Using the Bregman divergence, we can also describe MD as:

$$D\varphi(y^{t+1}) = D\varphi(x^t) - \eta_t g^t, \text{ with } g^t \in \partial f(x^t) \quad (5.3a)$$

$$x^{t+1} \in P_{C,\varphi}(y^{t+1}) = \underset{z \in C}{\operatorname{argmin}} D_\varphi(z, y^{t+1}) \quad (5.3b)$$

This performs gradient descent in certain "dual" space.

上面的式子可以用最优情况来观察一下：

$$D_\varphi(z, y^{t+1}) = \varphi(z) - \varphi(y^{t+1}) - \langle D\varphi(y^{t+1}), z - y^{t+1} \rangle$$

• (5.3b) 的最优情况是：

$$0 \in D\varphi(x^{t+1}) - D\varphi(y^{t+1}) + N_C(x^{t+1}) \quad \begin{array}{l} \text{通过因为} \\ \partial \varphi(x) = N_C(x) \\ (\text{subgradient}) \end{array}$$

最优时梯度为0 用向量替换↑ normal cone

$$(5.3a) = \varphi(x^{t+1}) - D\varphi(x^t) + \eta_t g^t + N_C(x^{t+1})$$

• (5.1)'s optimality condition gives :

这里用一下 5.1 的式子

$$x^{t+1} = \underset{x \in C}{\operatorname{argmin}} \left\{ \langle g^t, x - x^t \rangle + \frac{1}{\eta_t} D_\varphi(x, x^t) \right\}$$

with $g^t \in \partial f(x^t)$

$$0 \in \underbrace{g^t + \frac{1}{\eta_t} \{ D\varphi(x^{t+1}) - D\varphi(x^t) \}}_{N_C(x^{t+1})}$$

- Clearly, these two conditions are identical
 这两个只差一个系数 η_t , 因此这种形式与原形式是等价的

为了简化我们假设凸集 $C = \mathbb{R}^n$
 于是另一种形式是

$$x^{t+1} = D\varphi^*(D\varphi(x^t) - \eta g^t)$$

where $\varphi(x^*) := \sup_z \{ \langle z, x \rangle - \varphi(z) \}$ is the Fenchel-conjugate of φ

When $C = \mathbb{R}^n$, (5.3a)-(5.3b) simplifies to

$$x^{t+1} = y^{t+1} = \underbrace{(D\varphi)^{-1}(D\varphi(x^t) - \eta g^t)}_{\text{逆运算}}$$

因此, 可以得到这样的结论

$$(D\varphi)^{-1} = (D\varphi)^* \quad (5.5)$$

Proof of Claim (5.5):

Suppose $y = D\varphi(x)$

From the conjugate subgradient theorem, this is equivalent to:

$$\varphi(x) + \varphi^*(y) = \langle x, y \rangle$$

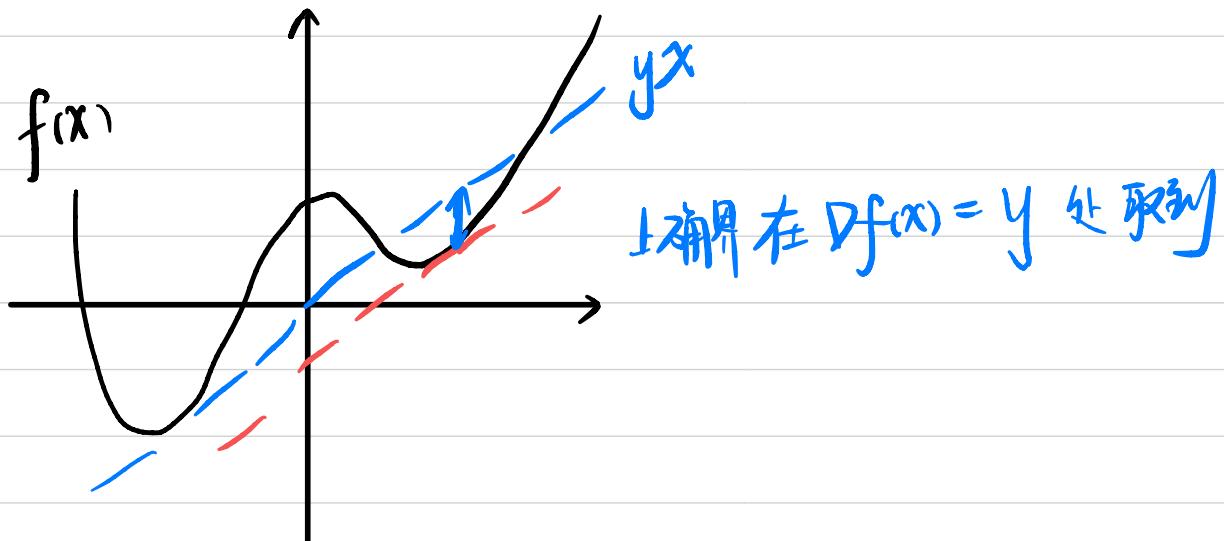
Conjugate subgradient theorem

Define (conjugate function):

For $f: \mathbb{R}^n \rightarrow \mathbb{R}$, $f^*: \mathbb{R}^n \rightarrow \mathbb{R}$ is the conjugate function of f

$$\text{iff } f^*(y) = \sup_{x \in \text{dom } f} (y^T x - f(x))$$

Intuition:



Theorem c (Conjugate Sub-gradient Theorem)

If f is closed and convex.

Then for a pair of vectors x, y , the followings are equivalent:

1. $x^T y = f(x) + f^*(y)$
2. $y \in \partial f(x)$
3. $x \in \partial f^*(y)$

We continue the proof:

$$(\nabla \varphi)^{-1} = (\nabla \varphi)^* \quad (5-5)$$

Proof of Claim (5.5) :

Suppose $y = \nabla \varphi(x)$

From the conjugate subgradient theorem, this is equivalent to:

$$\varphi(x) + \varphi^*(y) = \langle x, y \rangle$$

Since $\varphi^{**} = \varphi$, we further have

$$\varphi^*(y) + \varphi^{**}(x) = \langle x, y \rangle$$

which combined with the conjugate subgradient theorem, yields

$$x = \nabla \varphi^*(y)$$

Since $y = \nabla \varphi(x)$, we plug into:

$$x = \nabla \varphi^*(\nabla \varphi(x))$$

Hence,

$$(\nabla \varphi^*)^{-1}(x) = \nabla \varphi(x)$$

$$\text{Thus, } \nabla \varphi^* = \nabla \varphi^{-1}$$

Convergence analysis

Convex and Lipschitz problems

minimize $f(x)$

subject to $x \in C$

$\mathcal{S} = \{x - t\}$ 理想:

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 C . Suppose φ is ρ -strongly convex w.r.t. $\|\cdot\|$. Then

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

$$\|z\|_* = \sup \{ z^T x \mid \|x\| \leq 1 \}$$

我们可将带入一个直角一看

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

$$f^{\text{best},t} - f^{\text{opt}} \leq \frac{R + \frac{L_f^2}{2\rho} \sum_{k=0}^t \frac{2\rho R}{L_f^2 k}}{\sum_{k=0}^t \frac{\sqrt{2\rho R}}{L_f} \frac{1}{\sqrt{k}}}$$

$$= \frac{R + \sum_{k=0}^t \frac{R}{t}}{\sum_{k=0}^t \frac{1}{\sqrt{k}}} \cdot \frac{L_f}{\sqrt{2\rho R}} = O\left(\frac{L_f \sqrt{R}}{\sqrt{\rho}} \frac{\log t}{\sqrt{t}}\right)$$

Example: optimization over probability simplex

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

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

$$D_\varphi(x, x^0) = \varphi(x) - \varphi(x^0) - \langle D\varphi(x^0), x - x^0 \rangle$$

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

$$\because \|x\|_2^2 = x_1^2 + x_2^2 + \dots + x_n^2 \leq 1$$

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

$$\therefore R \leq \frac{1}{2}$$

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

再根据上面代入算出的结论, 以及 Theorem 5.3, 可得

$$f^{\text{best}} - f^{\text{opt}} \leq O(L_{f,2} \frac{\log t}{\epsilon})$$

(2) 再改變 $\phi(x) = -\sum_{i=1}^n x_i \log x_i$, which is $\underline{\text{strongly convex w.r.t. } \| \cdot \|_1}$, $\rho=1$

然後

(因為負熵的 D_p 等於 KL 散度)

$$\begin{aligned}\sup_{x \in \Delta} D_p(x, x^*) &= \sup_{x \in \Delta} \text{KL}(x \| x^*) \\ &= \sup_{x \in \Delta} \sum_{i=1}^n x_i \log x_i - \sum_{i=1}^n x_i \log \frac{1}{n}\end{aligned}$$

$$\sup_{x \in \Delta} -\sum_{i=1}^n x_i \log \frac{1}{n} = \log n \cdot \frac{1}{n} \cdot n = \log n$$

$$\sup_{x \in \Delta} \sum_{i=1}^n x_i \log x_i \leq 0 \quad \text{因為 } \log x_i \leq 0$$

因此：

$$\sup_{x \in \Delta} D_p(x, x^*) \leq \log n$$

FWL:

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

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

比較上面兩個不同 ℓ 的結果：

$$\text{Euclidean: } O(L_f, 2 \frac{\log t}{\epsilon})$$

$$KL: O(L_f, \sqrt{\log n} \cdot \frac{\log t}{\epsilon})$$

因為 $\|g\|_\infty = \max_i |g_i|$

$$\|g\|_2 = \sqrt{\sum_i g_i^2}$$

所以 $\|g\|_\infty \leq \|g\|_2 \leq \sqrt{n} \|g\|_\infty$

FFW. $\frac{1}{\sqrt{n}} \leq \frac{L_f, \infty}{L_f, 2} \leq 1$

所以 KL yields much better performance

速度更快

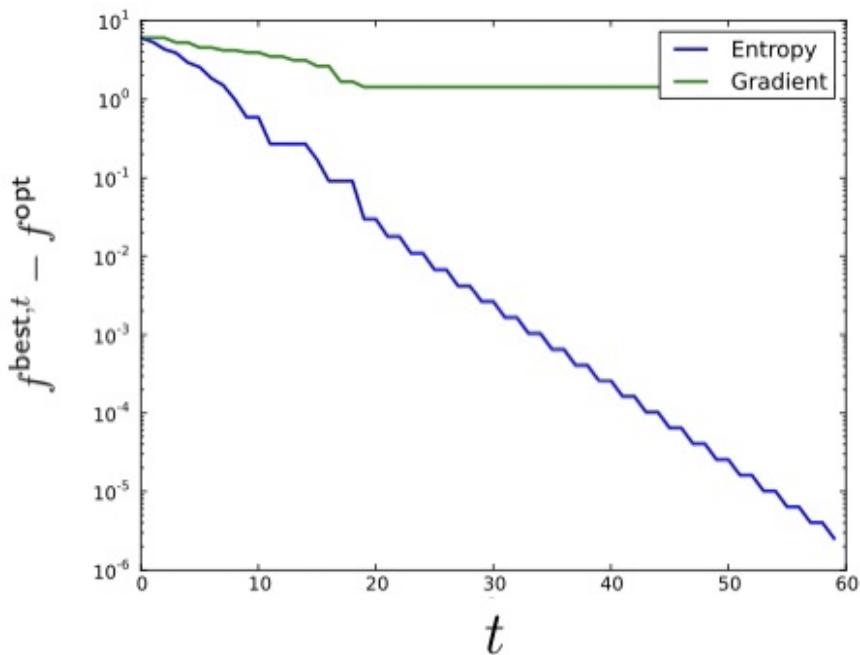
Numerical example : robust regression

$$\underset{x}{\text{minimize}} \quad f(x) = \sum_{i=1}^m |a_i^T x - b_i|$$

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

with $a_i \sim N(0, I_{n \times n})$, $b_i = \frac{a_{i,1} + a_{i,2}}{\sqrt{2}} + N(0, 10^{-2})$

$$m = 20 \quad n = 3000$$



Fundamental inequality for mirror descent

基本不等式

Lemma 5.4

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

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

Proof of Lemma 5.4:

$$\begin{aligned} f(\mathbf{x}^t) - f(\mathbf{x}^*) &\leq \langle g^t, \mathbf{x}^t - \mathbf{x}^* \rangle \quad (\text{Subgradient 等性质}) \\ &= \frac{1}{\eta_t} \langle \nabla \varphi(\mathbf{x}^t) - \nabla \varphi(\mathbf{y}^{t+1}), \mathbf{x}^t - \mathbf{x}^* \rangle \quad (\text{MD 的更新法}) \\ &= \frac{1}{\eta_t} \{ D_p(\mathbf{x}^*, \mathbf{x}^t) + D_p(\mathbf{x}^t, \mathbf{y}^{t+1}) - D_p(\mathbf{x}^*, \mathbf{y}^{t+1}) \} \end{aligned}$$

(three point lemma)

For every three points $\mathbf{x}, \mathbf{y}, \mathbf{z}$,

$$D_p(\mathbf{x}, \mathbf{z}) = D_p(\mathbf{x}, \mathbf{y}) + D_p(\mathbf{y}, \mathbf{z}) - \langle \nabla \varphi(\mathbf{z}) - \nabla \varphi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$$

$$\begin{aligned} &\leq \frac{1}{\eta_t} \{ D_p(\mathbf{x}^*, \mathbf{x}^t) + D_p(\mathbf{x}^t, \mathbf{y}^{t+1}) - D_p(\mathbf{x}^*, \mathbf{x}^{t+1}) \\ &\quad - D_p(\mathbf{x}^{t+1}, \mathbf{y}^{t+1}) \} \end{aligned}$$

(Pythagorean)

If $\mathbf{x}_{c,p} = P_{C,p}(\mathbf{x})$, then

$$D_p(\mathbf{z}, \mathbf{x}) \geq D_p(\mathbf{z}, \mathbf{x}_{c,p}) + D_p(\mathbf{x}_{c,p}, \mathbf{x}), \quad \forall \mathbf{z} \in C$$

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

反复極限 $\varphi(\cdot)$ 連續的， $\varphi(\cdot)$ bound 等於 L_f

即證：

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

證明。

$$\begin{aligned} & D_\varphi(x^t, y^{t+1}) - D_\varphi(x^{t+1}, y^{t+1}) \\ &= \varphi(x^t) - \varphi(x^{t+1}) - \langle \nabla \varphi(y^{t+1}), x^t - x^{t+1} \rangle \quad (\text{定義}) \\ &\leq \underbrace{\langle \nabla \varphi(x^t), x^t - x^{t+1} \rangle}_{\text{(strong convexity of } \varphi)} - \frac{\rho}{2} \|x^t - x^{t+1}\|^2 - \langle \nabla \varphi(y^{t+1}), x^t - x^{t+1} \rangle \quad (\Leftarrow \rho\text{-strong convexity}) \\ &= \langle \nabla \varphi(x^t) - \nabla \varphi(y^{t+1}), x^t - x^{t+1} \rangle - \frac{\rho}{2} \|x^t - x^{t+1}\|_2^2 \\ &= \eta_t \langle g^t, x^t - x^{t+1} \rangle - \frac{\rho}{2} \|x^t - x^{t+1}\|^2 \quad (\text{MD update rule}) \\ &\leq \eta_t L_f \|x^t - x^{t+1}\| - \frac{\rho}{2} \|x^t - x^{t+1}\|^2 \\ &\leq \frac{(\eta_t L_f)^2}{2\rho} \quad (\text{optimize quadratic function in } \|x^t - x^{t+1}\|) \end{aligned}$$

$\langle g^t, x^t - x^{t+1} \rangle \leq \|\eta_t\| \|x^t - x^{t+1}\|$
 $\varphi(x^t) \geq \varphi(x^t) + \langle \nabla \varphi(x^t), x^{t+1} - x^t \rangle + \frac{\rho}{2} \|x^{t+1} - x^t\|^2$

二次函数最大值在 $-\frac{b}{2a} = \frac{\eta_t L_f}{\rho}$ 处

MD update rule

$$D_\varphi(y^{t+1}) = D_\varphi(x^t) - \eta_t g^t, \text{ with } g^t \in \partial \varphi(x^t)$$

$$x^{t+1} \in P_{C,\varphi}(y^{t+1}) = \arg \min_{z \in C} D_\varphi(z, y^{t+1})$$

(5.3a)

(5.3b)

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

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

有了 Lemma 5.4 的，由 JL 3 点引理得证 by Theorem 5.3

From Lemma 5.4, one has

(Lemma 5.4) $\eta_k (f(x^k) - f^{\text{opt}}) \leq D_\varphi(x^*, x^k) - D_\varphi(x^*, x^{k+1}) + \frac{\eta_k^2 L_f^2}{2\rho}$

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

$\sum_{k=0}^t \eta_k (f(x^k) - f^{\text{opt}}) \leq D_\varphi(x^*, x^0) - D_\varphi(x^*, x^{t+1}) + \frac{L_f^2 \sum_{k=0}^t \eta_k^2}{2\rho}$

3-point lemma: $D_\varphi(x^*, x^0) \leq D(x^*, x^{t+1}) + D(x^{t+1}, x^0)$

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

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

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