

MATH466/MATH766

Math of machine learning

Appendix Convergence Analysis

References:

-

Todays contents: convergence analysis of

- gradient descent
- subgradient method
- proximal gradient
- stochastic gradient

Important concepts:

-

Recommend reading:

-

GD

$$\min_{x \in \mathbb{R}^n} f(x) \quad x^{(k+1)} = x^{(k)} - t \nabla f(x^{(k)})$$

Thm 1. Assume f is L -smooth, and $f(x) \geq f^*$ for any x

$$\text{then with } t = \frac{1}{L}, \min \left\{ \|\nabla f(x^{(k)})\|^2 \right\}_{k=0}^{K-1} \leq \frac{2L(f(x^{(0)}) - f^*)}{K}$$

Thm 2. In addition to Thm 1 assumptions,

if f is convex and $f(x) \geq f(x^*)$ for any x .

$$\text{then } f(x^{(k)}) - f(x^*) \leq \frac{L \|x^{(0)} - x^*\|^2}{2k}$$

Thm 3. In addition to Thm 1, 2 assumptions,

if f is μ -strongly convex.

$$\text{then } f(x^{(k)}) - f(x^*) \leq (1 - \frac{\mu}{L})^k (f(x^{(0)}) - f(x^*))$$

$$\text{and } \|x^{(k)} - x^*\|^2 \leq (1 - \frac{\mu}{L})^k \|x^{(0)} - x^*\|^2$$

Proof. 1. By Lipschitz smoothness,

$$\begin{aligned} f(x^{(k+1)}) &\leq f(x^{(k)}) - t \|\nabla f(x^{(k)})\|_2^2 + \frac{Lt^2}{2} \|\nabla f(x^{(k)})\|_2^2 \\ &= f(x^{(k)}) - t \left(1 - \frac{L}{2}\right) \|\nabla f(x^{(k)})\|_2^2 \end{aligned}$$

$$\text{Pick } t = \frac{1}{L}$$

$$f(x^{(k+1)}) \leq f(x^{(k)}) - \frac{1}{2L} \|\nabla f(x^{(k)})\|_2^2 \quad (*)$$

$$\Rightarrow \|\nabla f(x^{(k)})\|_2^2 \leq 2L(f(x^{(k)}) - f(x^{(k+1)}))$$

$$\begin{aligned} \Rightarrow \sum_{k=0}^{K-1} \|\nabla f(x^{(k)})\|_2^2 &\leq 2L(f(x^{(0)}) - f(x^{(K)})) \\ &\leq 2L(f(x^{(0)}) - f^*) < \infty \end{aligned}$$

$$\text{therefore } \min \left\{ \|\nabla f(x^{(k)})\|_2^2 \right\}_{k=0}^{K-1} \leq \frac{2L(f(x^{(0)}) - f^*)}{K}$$

2. From L-smooth, we have

$$f(x^{(k+1)}) \leq f(x^{(k)}) + \langle \nabla f(x^{(k)}), x^{(k+1)} - x^{(k)} \rangle + \frac{L}{2} \|x^{(k+1)} - x^{(k)}\|^2$$

By μ -convexity, we have

$$f(x^{(k)}) + \langle \nabla f(x^{(k)}), x^* - x^{(k)} \rangle + \frac{\mu}{2} \|x^{(k)} - x^*\|^2 \leq f(x^*)$$

Combine above two, we have

$$f(x^{(k+1)}) - f(x^*) \leq \langle \nabla f(x^{(k)}), x^{(k+1)} - x^* \rangle + \frac{L}{2} \|x^{(k+1)} - x^{(k)}\|^2 - \frac{\mu}{2} \|x^{(k)} - x^*\|^2$$

$$\text{Let } a := x^{(k)} - x^*, \quad b := x^{(k+1)} - x^*, \quad (\#)$$

$$\text{then } x^{(k+1)} - x^{(k)} = b - a, \quad \nabla f(x^{(k)}) = -\frac{1}{t}(b-a)$$

$$\text{RHS of } (\#) = -\frac{1}{t} \langle b-a, b \rangle + \frac{L}{2} \|b-a\|^2 - \frac{\mu}{2} \|a\|^2$$

$$\stackrel{t=\frac{1}{L}}{=} L \langle b-a, -b + \frac{1}{2}(b-a) \rangle$$

$$= -\frac{L}{2} \langle b-a, b+a \rangle$$

$$= \frac{L-\mu}{2} \|x^{(k)} - x^*\|^2 - \frac{L}{2} \|x^{(k+1)} - x^*\|^2 \quad (\Delta)$$

When $\mu=0$

$$f(x^{(k+1)}) - f(x^*) \leq \frac{L}{2} (\|x^{(k)} - x^*\|^2 - \|x^{(k+1)} - x^*\|^2)$$

$$\Rightarrow \sum_{k=0}^K (f(x^{(k)}) - f(x^*)) \leq \frac{L}{2} (\|x^{(0)} - x^*\|^2 - \|x^{(K+1)} - x^*\|^2)$$

$$\leq \frac{L}{2} \|x^{(0)} - x^*\|^2$$

Since $f(x^{(k+1)}) \leq f(x^{(k)})$

$$f(x^{(k)}) - f(x^*) \leq \frac{L \|x^{(0)} - x^*\|^2}{2K}$$

3. From 1, we have

$$f(x^{(k+1)}) - f(x^*) \leq f(x^{(k)}) - f(x^*) - t(1 - \frac{L}{2}) \|\nabla f(x^{(k)})\|^2 \quad (*)$$

Intuition: strong convexity.

\Rightarrow if $\|\nabla f(x)\|$ is small ($\|\nabla f(x)\|$ close to $\|\nabla f(x^*)\|$)
then x is close to x^* , $f(x)$ is close to $f(x^*)$

By strong convexity

$$h(y) := \underline{f(x^{(k)}) + \langle \nabla f(x^{(k)}), y - x^{(k)} \rangle + \frac{\mu}{2} \|y - x^{(k)}\|^2 \leq f(y)}, \forall y$$

quadratic in y

$$h(y) \geq h(x^{(k)} - \frac{1}{\mu} \nabla f(x^{(k)})) = f(x^{(k)}) - \frac{1}{2\mu} \|\nabla f(x^{(k)})\|^2$$

$$\Rightarrow -\|\nabla f(x^{(k)})\|^2 \leq 2\mu(f(x^*) - f(x^{(k)}))$$

$$\begin{aligned} \text{Plug in } (*) \quad f(x^{(k+1)}) - f(x^*) &\leq \left[1 - 2\mu t \left(1 - \frac{L}{2} \right) \right] (f(x^{(k)}) - f(x^*)) \\ &= (1 - 2\mu t + \mu L t^2) (f(x^{(k)}) - f(x^*)) \quad (***) \end{aligned}$$

$$\text{Let } C(t) := 1 - 2\mu t + \mu L t^2 = \mu L (t - \frac{1}{L})^2 + (1 - \frac{\mu}{L})$$

$$\text{When } t \in (0, \frac{2}{L}), \quad C(t) \in (0, 1),$$

$$\text{For largest shrinkage, take } t = \frac{1}{L}$$

$$f(x^{(k+1)}) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right) (f(x^{(k)}) - f(x^*))$$

$$\text{By induction} \quad f(x^{(k)}) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right)^k (f(x^{(0)}) - f(x^*))$$

From 2 (Δ)

$$0 \leq f(x^{(k+1)}) - f(x^*) \leq \frac{L - \mu}{2} \|x^{(k)} - x^*\|^2 - \frac{L}{2} \|x^{(k+1)} - x^*\|^2$$

$$\Rightarrow \|x^{(k+1)} - x^*\|^2 \leq \left(1 - \frac{\mu}{L}\right) \|x^{(k)} - x^*\|^2$$

$$\Rightarrow \|x^{(k)} - x^*\|^2 \leq \left(1 - \frac{\mu}{L}\right)^k \|x^{(0)} - x^*\|^2$$

SubGrad

$$\min_{x \in \mathbb{R}^n} f(x) \quad x^{(k+1)} = x^{(k)} - t_k g_k, \quad g_k \in \partial f(x^{(k)})$$

$$f(x_{\text{best}}^{(k)}) := \min_{k=0, \dots, K} f(x^{(k)})$$

Thm. Assume that f is L -Lipschitz continuous, i.e. $|f(x) - f(y)| \leq L \|x - y\|_2$

$$\text{Then } f(x_{\text{best}}^{(k)}) - f(x^*) \leq \frac{\|x^{(0)} - x^*\|^2 + L^2 \sum_{k=1}^K t_k^2}{2 \sum_{k=1}^K t_k}$$

Pf. By def of subgrad

$$f(x^*) \geq f(x^{(k)}) + \langle g_k, x^* - x^{(k)} \rangle$$

$$\begin{aligned} \Rightarrow f(x^{(k)}) - f(x^*) &\leq \langle g_k, x^* - x^{(k)} \rangle \\ &= \frac{1}{t_k} \langle x^{(k)} - x^{(k+1)}, x^{(k)} - x^* \rangle \end{aligned} \quad (*)$$

$$\text{Let } a := x^{(k)} - x^*, \quad b := x^{(k+1)} - x^*,$$

$$\text{then } b - a = x^{(k+1)} - x^{(k)} = -t_k g_k$$

$$\text{RHS of } (*) = \frac{1}{t_k} \langle a - b, \frac{a+b}{2} + \frac{a-b}{2} \rangle$$

$$= \frac{1}{2t_k} (\|a\|^2 - \|b\|^2) + \frac{t_k}{2} \|g_k\|^2$$

$$\Rightarrow 2t_k (f(x^{(k)}) - f(x^*)) \leq \|x^{(k)} - x^*\|^2 - \|x^{(k+1)} - x^*\|^2 + t_k^2 \|g_k\|^2$$

$$\Rightarrow \sum_{k=1}^K 2t_k (f(x^{(k)}) - f(x^*)) \leq \|x^{(0)} - x^*\|^2 + \sum_{k=1}^K t_k^2 \|g_k\|^2$$

$$\Rightarrow f(x_{\text{best}}^{(k)}) - f(x^*) \leq \frac{\|x^{(0)} - x^*\|^2 + \sum_{k=1}^K t_k^2 \|g_k\|^2}{2 \sum_{k=1}^K t_k}$$

$$\text{ProxGD} \quad \min_{x \in \mathbb{R}^n} f(x) = h(x) + r(x)$$

h differentiable,

r not necessarily differentiable

$$x^{(k+1)} = \text{prox}_{tr}(x^{(k)} - t \nabla h(x^{(k)}))$$

Thm. If h is convex, differentiable and L -smooth,

r is convex and $\text{prox}_{tr}(x)$ can be evaluated

then proximal gradient descent with fixed step size $t \leq \frac{1}{L}$

$$\text{satisfies } f(x^{(k)}) - f^* \leq \frac{\|x^{(0)} - x^*\|^2}{2tk}$$

Proof. Denote $G_t(x) := \frac{1}{t}(x - \text{prox}_{tr}(x - t \nabla h(x)))$

$$\text{Then } x^{(k+1)} = x^{(k)} - t G_t(x^{(k)})$$

(Central Lemma)

Claim: For any z ,

$$f(x^{(k+1)}) \leq f(z) + \langle G_t(x^{(k)}), x^{(k)} - z \rangle - (t - \frac{Lt^2}{2}) \|G_t(x^{(k)})\|^2$$

$$\text{Take } t = \frac{1}{L}$$

$$\text{Take } z = x^{(k)} \Rightarrow f(x^{(k+1)}) \leq f(x^{(k)}) - \frac{t}{2} \|G_t(x^{(k)})\|^2 \geq 0$$

$$\text{Take } z = x^*$$

make progress each iter

$$\Rightarrow f(x^{(k+1)}) - f(x^*) \leq \frac{1}{t} \langle x^{(k)} - x^{(k+1)}, x^{(k)} - x^* \rangle - \frac{1}{2t} \|x^{(k)} - x^{(k+1)}\|^2$$

similar strategy

$$f(x^{(k+1)}) - f(x^*) \leq \frac{1}{2t} (\|x^{(k)} - x^*\|^2 - \|x^{(k+1)} - x^*\|^2)$$

$$\Rightarrow f(x^{(k)}) - f(x^*) \leq \frac{\|x^{(0)} - x^*\|^2}{2tk}$$

Proof of Claim:

(Differentiable part)

By L-smooth, we have

$$h(x^{(k+1)}) \leq h(x^{(k)}) + \langle \nabla h(x^{(k)}), x^{(k+1)} - x^{(k)} \rangle + \frac{L}{2} \|x^{(k+1)} - x^{(k)}\|^2$$

By convexity of h , we have

$$h(x^{(k)}) + \langle \nabla h(x^{(k)}), z - x^{(k)} \rangle \leq h(z)$$

Combine above

$$\begin{aligned} \Rightarrow h(x^{(k+1)}) &\leq h(z) + \langle \nabla h(x^{(k)}), x^{(k)} - z \rangle \\ &\quad + \langle \nabla h(x^{(k)}), x^{(k+1)} - x^{(k)} \rangle + \frac{L}{2} \|x^{(k+1)} - x^{(k)}\|^2 \end{aligned} \quad (*)$$

(Non-differentiable part)

By convexity of r , for any $g \in \partial r(w^{(k+1)})$, z

$$r(x^{(k+1)}) + \langle g, z - x^{(k+1)} \rangle \leq r(z) \quad (**)$$

By update rule

$$x^{(k+1)} = \underset{x}{\operatorname{argmin}} \quad r(x) + \frac{1}{2t} \|x - (x^{(k)} + t \nabla h(x^{(k)}))\|^2$$

By optimality condition

$$0 \in \partial r(x^{(k+1)}) + \frac{1}{t} (x^{(k+1)} - x^{(k)} + t \nabla h(x^{(k)}))$$

$$\Rightarrow \frac{1}{t} (x^{(k)} - x^{(k+1)}) - \nabla h(x^{(k)}) \in \partial r(x^{(k+1)})$$

$$\Rightarrow G_t(x^{(k)}) - \nabla h(x^{(k)}) \in \partial r(x^{(k+1)})$$

Combine with (**)

$$r(x^{(k+1)}) \leq r(z) + \langle \nabla h(x^{(k)}) - G_t(x^{(k)}), z - x^{(k+1)} \rangle \quad . \forall z \quad (***)$$

(*) + (***)

$$f(x^{(k+1)}) \leq f(z) + \langle G_t(x^{(k)}), x^{(k+1)} - z \rangle + \frac{L}{2} \|x^{(k+1)} - x^{(k)}\|^2$$

$$= f(z) + \langle G_t(x^{(k)}), x^{(k)} - z \rangle - (t - \frac{Lt^2}{2}) \|G_t(x^{(k)})\|^2$$

SGD

$$\min_{\mathbf{w}} f(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}) \quad \text{sample from } \{1, 2, \dots, n\}$$

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \alpha_k \nabla f_{i_k}(\mathbf{w}^k)$$

Thm. If f is L -smooth, $f(\mathbf{w}) \geq f^*$ for any \mathbf{w}

$$\mathbb{E}_i [\|\nabla f_i(\mathbf{w})\|^2] \leq G^2$$

$$\text{Then } \min \left\{ \mathbb{E} [\|\nabla f(\mathbf{w}^k)\|^2] \right\}_{k=0}^{K-1} \leq \frac{f(\mathbf{w}^0) - f^* + \frac{\alpha_k^2 L}{2} \sum_{k=0}^{K-1} t_k^2}{\sum_{k=0}^{K-1} t_k}$$

Proof. By L -smooth

$$f(\mathbf{w}^{k+1}) \leq f(\mathbf{w}^k) - t_k \nabla f(\mathbf{w}^k)^T \nabla f_{i_k}(\mathbf{w}^k) + \frac{t_k^2 L}{2} \|\nabla f_{i_k}(\mathbf{w}^k)\|_2^2$$

may not be negative

But "on average" (in expectation), we may make progress

$$I_{k-1} := (i_1, i_2, \dots, i_{k-1})$$

$$\begin{aligned} \mathbb{E}_{I_k} [f(\mathbf{w}^{k+1})] &\leq \mathbb{E}_{I_k} [f(\mathbf{w}^k) - t_k \nabla f(\mathbf{w}^k)^T \nabla f_{i_k}(\mathbf{w}^k) + \frac{t_k^2 L}{2} \|\nabla f_{i_k}(\mathbf{w}^k)\|_2^2] \\ &= \mathbb{E}_{I_{k-1}} [f(\mathbf{w}^k)] - t_k \mathbb{E}_{I_{k-1}} [\|\nabla f(\mathbf{w}^k)\|^2] + \frac{t_k^2 L}{2} \mathbb{E}_{i_k} [\|\nabla f_{i_k}(\mathbf{w}^k)\|_2^2] \\ &\leq \mathbb{E} [f(\mathbf{w}^k)] - t_k \mathbb{E} [\|\nabla f(\mathbf{w}^k)\|^2] + \frac{t_k^2 G^2 L}{2} \end{aligned}$$

$$\Rightarrow \sum_{k=0}^{K-1} t_k \mathbb{E} [\|\nabla f(\mathbf{w}^k)\|^2] \leq f(\mathbf{w}^0) - f^* + \frac{G^2 L}{2} \sum_{k=0}^{K-1} t_k^2$$

$$\Rightarrow \min \left\{ \mathbb{E} [\|\nabla f(\mathbf{w}^k)\|^2] \right\}_{k=0}^{K-1} \leq \frac{f(\mathbf{w}^0) - f^* + \frac{G^2 L}{2} \sum_{k=0}^{K-1} t_k^2}{\sum_{k=0}^{K-1} t_k}$$