

Lecture 15. Newton's method cont'd; Proximal gradient descent

Recall the convergence analysis of univariate functions.

$$|x_{k+1} - x^*| \leq C |x_k - x^*|^2 \text{ by Taylor expansion if } |f''(x)| \leq M.$$

$$\text{also proved by a slightly loose condition } |f''(x) - f''(y)| \leq M|x-y|.$$

To define M -Lipschitz for $\nabla^2 f(x)$, where $f: \mathbb{R}^n \rightarrow \mathbb{R}$, define norm first.

$$\text{Operator norm: } \|Q\|_{a,b} \triangleq \max_{x \neq 0} \frac{\|Qx\|_b}{\|x\|_a} \Rightarrow \|Qx\|_b \leq \|Q\|_{a,b} \|x\|_a$$

In particular, if $a=b=2$. $\|\cdot\|_{a,b}$ is called spectral norm. $\|\cdot\|_2$.

Proposition. $\|Q\|_2 = (\lambda_{\max}(Q^T Q))^{1/2}$ since $\|Qx\|^2 = x^T Q^T Q x \leq \lambda_{\max} \|x\|^2$.

$\nabla^2 f(x)$ is M -Lipschitz if $\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq M \|x-y\|_2$. $\forall x, y$.

Theorem. If f is m -strongly convex. $\nabla^2 f$ is M -Lipschitz. $\{x_k\}$ produced

by Newton's method. then $\|x_{k+1} - x^*\| \leq \frac{M}{2m} \|x_k - x^*\|^2$.

Remark: let $y_k = \frac{M}{2m} \|x_k - x^*\|$. then $y_{k+1} \leq y_k^2 \Rightarrow y_k \leq y_0^{2^k}$.

Proof. $\|x_{k+1} - x^*\|_2 \leq \|(\nabla^2 f(x_k))^{-1}\|_2 \|\nabla^2 f(x_k)(x_k - x^*) - (\nabla f(x_k) - \nabla f(x^*))\|_2$.

$$= \|(\nabla^2 f(x_k))^{-1}\|_2 \left\| \int_0^1 (\nabla^2 f(x_k) - \nabla^2 f(x^* + t(x_k - x^*))) (x_k - x^*) dt \right\|_2.$$

$$\leq \frac{1}{m} \|x_k - x^*\|_2 \cdot \int_0^1 \|\nabla^2 f(x_k) - \nabla^2 f(x^* + t(x_k - x^*))\|_2 dt.$$

$$\leq \frac{M}{m} \|x_k - x^*\|_2 \cdot \int_0^1 (1-t) \|x_k - x^*\|_2 dt = \frac{M}{2m} \|x_k - x^*\|_2^2. \quad \square$$

In summary. descent method $x_{k+1} = x_k - t_k \cdot d_k$. $x_k \rightarrow x^*$

direction: gradient descent $x_{k+1} = x_k - t_k \nabla f(x_k)$.

Newton's method $x_{k+1} = x_k - t_k (\nabla^2 f(x_k))^{-1} \nabla f(x_k)$.

step size: fixed step size / exact line search / backtracking

Convergence: by Lyapunov's stability theorem. If $f(x_{k+1}) < f(x_k)$.

order and rate : $|f(x_{k+1}) - f(x^*)| \leq \alpha |f(x_k) - f(x^*)|$ for some $\alpha < 1$.

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - x^*|}{|x_n - x^*|^\rho} = \mu.$$

linear convergence for gradient descent.

$$\|x_{k+1} - x^*\| \leq c \|x_k - x^*\|^2 \text{ for some } c < \infty$$

ρ : order of convergence

μ : rate of convergence. quadratic convergence for Newton's method.

condition : gradient descent

Newton's method.

L -smooth \Rightarrow convergence .

$\nabla^2 f(x)$ M-Lipschitz

m -Strongly convex. \Rightarrow rate.

m-strongly convex

How to deal with non differentiable functions: Proximal gradient descent.

consider gradient descent : $x_{k+1} = x_k - t \nabla f(x)$

'if f is L -smooth, and

$$x_{k+1} = \arg \min_y \quad f(x_k) + \nabla f(x_k)^T (y - x_k) + \frac{1}{2t} \|y - x_k\|^2.$$

$t \leq 1/L$. RHS gives an upper bound of f .

more sophisticated. $x_{k+1} = \arg \min_y f(x_k) + \nabla f(x_k)^T(y - x_k) + \frac{1}{2} (y - x_k)^T \nabla^2 f(x_k) (y - x_k)$.

Now assume f is not differentiable. suppose $f(x) = g(x) + h(x)$.

where $g(x)$ is convex and differentiable. and $h(x)$ is convex.

$$\begin{aligned} \text{let } x_{k+1} &= \underset{y}{\operatorname{argmin}} \quad g(x_k) + \nabla g(x_k)^T(y - x_k) + \frac{1}{2t} \|y - x_k\|^2 + h(y). \\ &= \underset{y}{\operatorname{argmin}} \quad \frac{1}{2t} \|y - (x_k + t \nabla g(x_k))\|^2 + h(y). \end{aligned}$$

$\frac{1}{2t} \|y - (x_k + t \nabla g(x_k))\|^2$: closed to update for g . $h(y)$: make h small.

$$\text{Proximal mapping: } \text{prox}_{h,t}(x) = \underset{y}{\operatorname{argmin}} \quad \frac{1}{2t} \|x - y\|^2 + h(y)$$

$$\text{Proximal gradient descent: } x_{k+1} = \text{prox}_{h,t}(x_k - t \cdot \nabla g(x_k)).$$

$$x_{k+1} = x_k - t_k G_{t_k}(x_k) \text{ where } G_t(x) = \frac{1}{t}(x - \text{prox}_{h,t}(x - t \cdot \nabla g(x))).$$

Key point is that $\text{prox}_{h,t}(x)$ has closed form for many important h .

Lasso: least absolute shrinkage and selection operator. (to reduce overfitting).

$$f(\beta) = \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1, \quad h(\beta) = \lambda \|\beta\|_1, \quad (\lambda > 0).$$

$$\text{Proximal mapping } \text{prox}_{h,t}(\beta) = \underset{y}{\operatorname{argmin}} \quad \frac{1}{2t} \|\beta - y\|_2^2 + \lambda \|y\|_1$$

$$\text{prox}_{h,t}(\beta) = \underset{y}{\operatorname{argmin}} \quad \sum_i \left(\frac{1}{2t} (\beta_i - y_i)^2 + \lambda |y_i| \right).$$

$$\Rightarrow [\text{prox}_{h,t}(\beta)]_i = \underset{y_i}{\operatorname{argmin}} \quad \frac{1}{2t} (\beta_i - y_i)^2 + \lambda |y_i| = \begin{cases} \beta_i - \lambda t & \text{if } \beta_i > \lambda t \\ 0 & \text{if } |\beta_i| \leq \lambda t \\ \beta_i + \lambda t & \text{if } \beta_i < -\lambda t \end{cases}$$

$$\Rightarrow \text{prox}_{h,t}(\beta) = S_{\lambda t}(\beta) \text{. the soft thresholding operator } \nearrow$$

$$\begin{aligned} \text{proximal update } \beta_{k+1} &= \text{prox}_{h,t}(\beta_k - t \nabla g(\beta_k)) = S_{\lambda t}(\beta_k + t X^T(y - X\beta_k)) \\ &= -X^T(y - X\beta_k) \end{aligned}$$

Remark: it is called the iterative soft-thresholding algorithm (ISTA).

Subgradient: note $f(y) \geq f(x) + \nabla f(x)^T(y-x)$. how about non-differentiable f ?

d_x is a subgradient. denoted by $d_x \in \partial f(x)$, if $f(y) \geq f(x) + d_x^T(y-x)$.

$x_{k+1} = x_k - t \cdot G_t(x_k)$ $G_t(x_k)$ is gradient or subgradient of $f=g+h$? \times

$$w = \text{prox}_{h,t} = \arg \min_y \frac{1}{2t} \|y-x\|^2 + h(y) \triangleq \tilde{h}(y). \quad 0 \in \partial \tilde{h}(w).$$

$$\frac{1}{t}(w-x) \in \partial(\frac{1}{2t}\|y-x\|^2)(w). \quad \text{conj. } \frac{1}{t}(x-w) \in \partial h(w).$$

$$\text{note } \tilde{h}(y) \geq \tilde{h}(w) \iff h(y) \geq h(w) + \underbrace{\frac{1}{2t}(2x-w-y)^T(y-w)}_{\stackrel{>0}{\leq}} \stackrel{>0}{\leq} 2(x-w)^T(y-w).$$

$$\text{assume } h(z) = h(w) + \frac{1}{t}(x-w)^T(z-w) - \delta \|z-w\|.$$

$$\text{by convexity. } h(y) \leq h(w) + \frac{1}{t}(x-w)^T(y-w) - \delta \|y-w\|. \quad y = \theta z + \bar{\theta} w.$$

$$\text{however. } \frac{1}{2t}\|y-x\|^2 = \frac{1}{2t}\|w-x\|^2 + \frac{1}{t}(w-x)(y-w) + \frac{1}{2t}\|y-w\|^2.$$

$$\Rightarrow h(y) + \frac{1}{2t}\|y-x\|^2 < h(w) + \frac{1}{2t}\|w-x\|^2 \text{ for sufficiently small } \|y-w\|. \quad \text{contradiction!}$$

In general. we have $\partial(f_1+f_2) = \partial f_1 + \partial f_2$.

$$\text{Firmly nonexpansive: } (\text{prox}_h(x) - \text{prox}_h(y))^T(x-y) \geq \|\text{prox}_h(x) - \text{prox}_h(y)\|^2$$

let $u = \text{prox}_h(x)$. $v = \text{prox}_h(y)$. then $x-u \in \partial h(u)$. $y-v \in \partial h(v)$.

$$\Rightarrow h(u) \geq h(w) + (x-u)^T(v-u). \text{ and } h(u) \geq h(w) + (y-v)^T(u-v).$$

$$\Rightarrow (x-u)^T(v-u) + (y-v)^T(u-v) \leq 0 \Rightarrow (x-u-y+v)^T(u-v) \geq 0.$$

Nonexpansive : $\|\text{prox}_h(x) - \text{prox}_h(y)\| \leq \|x - y\|$ by Cauchy-Schwarz.

Convergence analysis : assume g is L -smooth and m -strongly convex.

$$\text{Recall } \frac{m}{2} \|x - y\|^2 \leq g(y) - g(x) - \nabla g(x)^T (y - x) \leq \frac{L}{2} \|x - y\|^2. \quad y \leftarrow x_{k+1}$$

$$\frac{m}{2} t^2 \|G_t(x_k)\|^2 \leq g(x_{k+1}) - g(x_k) + t \nabla g(x_k)^T G_t(x_k) \leq \frac{L}{2} t^2 \|G_t(x_k)\|^2$$

$$\Rightarrow g(x_{k+1}) \leq g(x_k) - t \nabla g(x_k)^T G_t(x_k) + \frac{t}{2} \|G_t(x_k)\|^2 \quad (\text{select } t \leq 1/L)$$

$$\stackrel{y}{\Rightarrow} f(x_{k+1}) \leq g(y) - \nabla g(x_k)^T (y - x_k) - \frac{m}{2} \|x_k - y\|^2 + h(x_{k+1})$$

$$(\text{note } x_{k+1} = x_k - t G_t(x_k) = \text{prox}_{h,t}(x_k - t \nabla g(x_k)) \Rightarrow G_t(x_k) - \nabla g(x_k) \in \partial h(x_{k+1}))$$

$$\Rightarrow f(x_{k+1}) \leq g(y) + h(y) - G_t(x_k)^T (y - x_{k+1}) - \frac{m}{2} \|x_k - y\|^2 + \nabla g(x_k)^T (x_k - x_{k+1}).$$

$$f(y) = g(y) + h(y) - G_t(x_k)^T (y - x_k) - \frac{m}{2} \|x_k - y\|^2 - \frac{t}{2} \|G_t(x_k)\|^2 = t G_t(x_k).$$

$$\text{let } y = x_k \Rightarrow f(x_{k+1}) \leq f(x_k) - \frac{t}{2} \|G_t(x_k)\|^2. \quad (\text{descent method})$$

$$\text{let } y = x^* \Rightarrow f(x_{k+1}) - f(x^*) \leq G_t(x_k)^T (x_k - x^*) - \frac{m}{2} \|x_k - x^*\|^2 - \frac{t}{2} \|G_t(x_k)\|^2$$

$$= G_t(x_k)^T (x_k - x^* - \frac{t}{2} G_t(x_k)) - \frac{m}{2} \|x_k - x^*\|^2.$$

$$= \frac{1}{2t} (x_k - x_{k+1})^T ((x_k - x^*) + (x_{k+1} - x^*)) - \frac{m}{2} \|x_k - x^*\|^2.$$

$$\leq \frac{1}{2t} (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2) - \frac{m}{2} \|x_k - x^*\|^2.$$

$$\Rightarrow \sum_{k=0}^T (f(x_{k+1}) - f(x^*)) \leq \frac{1}{2t} \|x_0 - x^*\|^2 \Rightarrow f(x_T) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2tT}.$$

$$\text{Furthermore, } m > 0 \Rightarrow \|x_{k+1} - x^*\|^2 \leq (1 - mt) \|x_k - x^*\|^2. \quad \text{fix } t = 1/L$$

Require L -smooth. if L is unknown? Backtracking line search.

choose $0 < \alpha, \beta < 1$, and initial $\hat{t} > 0$. often choose $\alpha = 1/2$.

Recall in analysis we need $g(y) - g(x) - \nabla g(x)^T(y-x) \leq \frac{L}{2} \|x-y\|^2$

while $g(x_k - t_k G_{t_k}(x_k)) > g(x_k) - t \nabla g(x_k)^T G_{t_k}(x_k) + \frac{t_k}{2} \|G_{t_k}(x_k)\|^2$ $t_k = \beta t_k$

Convergence analysis for backtracking line search. note $t_k \geq \min\{\hat{t}, \beta/L\}$

$$f(x_{k+1}) - f(x^*) \leq \frac{1}{2t_k} ((1-mt_k) \|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2). \quad t_{\min} \triangleq$$

$$\Rightarrow \sum_{k=0}^{T-1} t_k (f(x_{k+1}) - f(x^*)) \leq \frac{1}{2} (\|x_0 - x^*\|^2 - \|x_T - x^*\|^2)$$

$$\Rightarrow (\sum t_k) (f(x_T) - f(x^*)) \leq \frac{1}{2} \|x_0 - x^*\|^2 \Rightarrow f(x_T) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2t_{\min}}$$

Again, if m -strongly convex for $m > 0$. $\|x_{k+1} - x^*\|^2 \leq (1-mt_{\min}) \|x_k - x^*\|^2$.