

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

Recall the convergence analysis of univariate functions.

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

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.

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

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.

ρ : order of convergence

μ : rate of convergence.

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

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

} convergence and rate.

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

$$x_{k+1} = \operatorname{argmin}_y f(x_k) + \nabla f(x_k)^T (y - x_k) + \frac{1}{2t} \|y - x_k\|^2$$

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

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

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.

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

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

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

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

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

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

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

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

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

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

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

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

