

-5- Nonlinear Optimization in Machine Learning -1-

Lecture 2 Line search methods

$$\min_{x \in \mathbb{R}^n} f(x)$$

$$\{x^k\} \quad f(x^{k+1}) < f(x^k) \quad k=0, 1, 2, \dots$$

"descent direction" $d \in \mathbb{R}^n$

$$f(x + td) < f(x) \text{ for all } t > 0 \text{ sufficient}$$

Say f is continuously differentiable

$$f(x + td) = f(x) + t \nabla f(x + \gamma td)^T d, \forall t \in (0, 1)$$

if $d^T \nabla f(x) < 0$ then d is a descent direction

"steepest descent"

$$\inf_{\|d\|=1} d^T \nabla f(x) = -\|\nabla f(x)\|$$

achieved when $d = -\frac{\nabla f(x)}{\|\nabla f(x)\|}$

"Steepest descent method" at minimizing residual

$$x^{k+1} = x^k - \alpha_k \nabla f(x^k), k=0,1,2,\dots$$

$\alpha_k > 0$ step length

"fixed point iteration"

$$x^{k+1} = \Phi(x^k) \quad k=0,1,2,\dots$$

$$\Phi(x) = x - \alpha \nabla f(x) \quad \text{for } \alpha > 0$$

$$\|\Phi(x) - \Phi(z)\| = \|(x-z) - \alpha(\nabla f(x) - \nabla f(z))\|$$

$$= \left\| (x-z) - \alpha \int_0^1 \nabla^2 f(z + \gamma(x-z))(x-z) d\gamma \right\|$$

$$= \left\| \left(I - \alpha \int_0^1 \nabla^2 f(z + \gamma(x-z)) d\gamma \right) (x-z) \right\|$$

$$\lambda \left[I - \alpha \int_0^1 \nabla^2 f(z + \gamma(x-z)) d\gamma \right] \in [I - \alpha L, I - \alpha m]$$

If x^* is s.t. $\nabla f(x^*) = 0$

and $\|\phi(x) - \phi(z)\| \leq \beta \|x - z\| \quad (\beta \in [0, 1])$

then $\|x^{k+1} - x^*\| = \|\phi(x^k) - \phi(x^*)\| \geq$

$$\leq \beta \|x^k - x^*\|$$

$$\leq \dots \leq \beta^{k+1} \|x_0 - x^*\|$$

"linear convergence rate"

In order $\|x^k - x^*\| \leq \varepsilon$ need $T \geq \frac{\log\left(\frac{\|x^0 - x^*\|}{\varepsilon}\right)}{\log \beta}$

We want $-\beta \leq 1 - \alpha L \leq 1 - \alpha m \leq \beta$

so $\alpha \in \left[\frac{1-\beta}{m}, \frac{1+\beta}{L}\right]$

When $1 - \alpha L = 1 - \alpha m \Rightarrow \beta = \frac{L-m}{L+m}, \alpha = \frac{2}{L+m}$

"steepest descent method"

$$\begin{aligned}
 f(x + \alpha d) &= f(x) + \alpha \nabla f(x)^T d + \alpha \int_0^1 [(\nabla f(x + \alpha d)) - (\nabla f(x))] d \gamma \\
 &\leq f(x) + \alpha \nabla f(x)^T d + \alpha \int_0^1 \|(\nabla f(x + \alpha d)) - (\nabla f(x))\| \|d\| d \gamma \\
 &\leq f(x) + \alpha \nabla f(x)^T d + \alpha \frac{\frac{L}{2} \|d\|^2}{2}
 \end{aligned}$$

$$x = x^k \quad d = -\nabla f(x^k)$$

$$\text{take } \alpha = \frac{1}{L}, \quad x^{k+1} = x^k - \frac{1}{L} \nabla f(x^k), \quad k=0,1,2.$$

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

We derive various convergence rate estimates from

(*)

General Case

$$f(x) \geq \bar{f} \quad \text{for all } x$$

$$\sum_{k=0}^{T-1} \|\nabla f(x^k)\|^2 \leq 2L \cdot \left(\sum_{k=0}^{T-1} [f(x^k) - f(x^{k+1})] \right)$$

$$= 2L \cdot [f(x^0) - f(x^T)]$$

$$\frac{\|x^T - x^0\|}{T} \leq (2L[f(x^0) - \bar{f}])$$

so $\lim_{T \rightarrow \infty} \|\nabla f(x_T)\| = 0$

$$\min_{0 \leq k \leq T-1} \|\nabla f(x^k)\| \leq \sqrt{\frac{2L[f(x^0) - \bar{f}]}{T}}$$

• Convex Case

$$f(x^*) \geq f(x^k) + \nabla f(x^k)^T (x^* - x^k)$$

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

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

$$= f(x^*) + \frac{L}{2} \left(\|x^k - x^*\|^2 - \|x^k - x^* - \frac{1}{L} \nabla f(x^k)\|^2 \right)$$

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

$$\text{So } \sum_{k=0}^{T-1} (f(x^{k+1}) - f^*) \leq \frac{L}{2} \sum_{k=0}^{T-1} \left(\|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2 \right) \quad \text{--- 6-}$$

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

as $f(x^k) \downarrow$ we get

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

- Strongly convex case

$$f(x) + (\nabla f(x))^T (z-x) + \frac{L}{2} \|z-x\|^2 \geq f(z) \geq f(x) + (\nabla f(x))^T (z-x) + \frac{m}{2} \|z-x\|^2$$

"sandwiching"

$$f_\mu(x) = f(x) + \mu \|x\|^2 \quad \text{convex} \rightarrow \text{strong convex}$$

Look at

$$f(z) \geq f(x) + (\nabla f(x))^T (z-x) + \frac{m}{2} \|z-x\|^2$$

$$\min_z f(z) \geq \min_z \left(f(x) + (\nabla f(x))^T (z-x) + \frac{m}{2} \|z-x\|^2 \right)$$

(at $z = x - \frac{\nabla f(x)}{m}$)

$$\Rightarrow f(x^*) \geq f(x) - \frac{1}{2m} \|\nabla f(x)\|^2$$

-8-
with some 3 steps. When it's standard case

i.e. $\|\nabla f(x)\|^2 \geq 2m[f(x) - f(x^*)]$

"linear convergence rate"

$$f(x^{k+1}) = f\left(x^k - \frac{1}{L} \nabla f(x^k)\right)$$

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

$$\leq f(x^k) - \frac{m}{L} (f(x^k) - f^*)$$

i.e. $f(x^{k+1}) - f(x^*) \leq \left(1 - \frac{m}{L}\right) (f(x^k) - f(x^*))$

$$\Rightarrow f(x^T) - f^* \leq \left(1 - \frac{m}{L}\right) (f(x^0) - f^*)$$

Iterate k s.t. $0 \leq k \leq T \quad \|\nabla f(x^k)\| \leq \varepsilon$

General $T \geq \frac{2L(f(x^0) - f^*)}{\varepsilon^2}$ sublinear

Weakly convex $T \leq \frac{f(x^0) - f^*}{\varepsilon}$ sublinear

Strongly convex $k \geq \frac{L}{m} \log\left(\frac{f(x^0) - f^*}{\varepsilon}\right)$ linear

-8-

Descent Methods with variable stepsize & search direction

$$x^{k+1} = x^k + \alpha_k d^k \quad k=0, 1, 2, \dots$$

$\alpha_k > 0$ d^k — search direction

$$-(d^k)^T \nabla f(x^k) \geq \bar{\epsilon} \|\nabla f(x^k)\| \cdot \|d^k\|$$

$$\gamma_1 \|\nabla f(x^k)\| \leq \|d^k\| \leq \gamma_2 \|\nabla f(x^k)\|$$

where $\bar{\epsilon}, \gamma_1, \gamma_2 > 0$

If $d^k = -\nabla f(x^k)$, then $\bar{\epsilon} = \gamma_1 = \gamma_2 = 1$

$$f(x^{k+1}) = f(x^k + \alpha d^k)$$

$$= f(x^k) + \alpha \nabla f(x^k)^T d^k + \alpha \frac{1}{2} \nabla f(x^k + \alpha d^k)^T \nabla f(x^k + \alpha d^k)$$

$$\leq f(x^k) + \alpha (\nabla f(x^k))^T d^k + \alpha \frac{1}{2} \|\nabla f(x^k + \alpha d^k) - \nabla f(x^k)\| \cdot \|d^k\| \alpha \gamma$$

$$\leq f(x^k) - \alpha \bar{\epsilon} \|\nabla f(x^k)\| \cdot \|d^k\| + \frac{\alpha^2}{2} L \|d^k\|^2$$

$$\leq f(x^k) - \alpha \left(\bar{\epsilon} - \alpha \frac{L}{2} \gamma_2 \right) \|\nabla f(x^k)\| \cdot \|d^k\|$$

If $\alpha \in (0, \frac{2\bar{\epsilon}}{L\gamma_2})$ then $f(x^{k+1}) < f(x^k)$

① $d^k = -S^k \nabla f(x^k)$ S^k symmetric positive definite
 $\lambda(S^k) \in [\gamma_1, \gamma_2]$

② Gauss-Southwell: $d^k = -[\nabla f(x^k)]_{i_k}$

$$i_k = \arg \min_i |[\nabla f(x^k)]_i|$$

③ Stochastic Coordinate Descent

$d^k = -[\nabla f(x^k)]_{i_k}$ $i_k \in \{1, 2, \dots, n\}$ chosen uniformly

④ Stochastic Gradient methods

⑤ Exact line search $\min_{\alpha > 0} f(x^k + \alpha d^k)$

⑥ Approximate Line search

"Weak Wolfe Conditions"

$$f(x^k + \alpha d^k) \leq f(x^k) + c_1 \alpha \nabla f(x^k)^T d^k \quad \text{"sufficient decrease"}$$

$$\nabla f(x^k + \alpha d^k)^T d^k \geq c_2 \nabla f(x^k)^T d^k$$

⑦ Backtracking line search

$$\bar{x}, \beta\bar{x}, \beta^2\bar{x}, \beta^3\bar{x}, \dots \left(\frac{\bar{x}}{3}, 0 \right) \rightarrow x$$

How to obtain convergence?

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

for some $C > 0$

Approximate Second-Order Necessary Points

$$\|\nabla f(x)\| \leq \varepsilon_g, \lambda_{\min}(\nabla^2 f(x)) \geq -\varepsilon_H \quad (**)$$

(i) If $\|\nabla f(x^k)\| > \varepsilon_g$ take the steepest descent

$$x^{k+1} = x^k - \frac{1}{L} \nabla f(x^k)$$

(ii). Otherwise let λ_k be the minimum eigenvalue of $\nabla^2 f(x^k)$. If $\lambda_k \leq -\varepsilon_H$ choose p^k to be the eigenvector of the most negative eigenvalue of $\nabla^2 f(x^k)$

$$\|p^k\| = 1, (p^k)^T \nabla f(x^k) \leq 0$$

$$\text{Set } x^{k+1} = x^k + \alpha_k p^k, \alpha_k = \frac{2/\lambda_k}{M}$$

(iii) If both (i) & (ii) are false then (**) is true

For (i)

$$\begin{aligned} f(x^{k+1}) &\leq f(x^k) - \frac{1}{2L} \|Df(x^k)\|^2 \\ &\leq f(x^k) - \frac{\epsilon_g^2}{2L} \end{aligned}$$

For (ii)

$$\begin{aligned} f(x^{k+1}) &\leq f(x^k) + \alpha_k Df(x^k)^T p^k + \frac{1}{2} \alpha_k^2 (p^k)^T D^2 f(x^k) p^k \\ &\quad + \frac{1}{6} M \alpha_k^3 \|p^k\|^3 \\ &\leq f(x^k) - \frac{1}{2} \left(\frac{2|\lambda_k|}{M} \right)^2 |\lambda_k| + \frac{1}{6} M \left(\frac{2|\lambda_k|}{M} \right)^3 \\ &= f(x^k) - \frac{2}{3} \frac{|\lambda_k|^3}{M^2} = f(x^k) - \frac{2}{3} \frac{\epsilon_H^3}{M^2} \end{aligned}$$

Here $\|D^2 f(x) - D^2 f(y)\| \leq M \|x - y\|$.

$$f(x+p) \leq f(x) + Df(x)^T p + \frac{1}{2} p^T D^2 f(x) p + \frac{1}{6} M \|p\|^3$$