

Problems with Newton's method

- 1) Calculating & storing Hessian is expensive for large problems
- 2) Solving the Newton systems can be expensive
- 3) We only get local convergence: initial point needs to be suff close to true solution!

1,2: Quasi-Newton methods

3: Globalization strategies (line search and trust region).

Line search methods

Def (descent direction) Let f cont diffble, $p \in \mathbb{R}^n$ is a descent dir if

$$\nabla f(x)^T p < 0$$

For descent directions: $\exists \epsilon > 0 \quad \forall t \in [0, \epsilon] \quad f(x+tp) < f(x)$

i.e. we can improve objective function if we look in dir p .

Examples:

$$p = -\nabla f(x)$$

$$p = -H^{-1} \nabla f(x) \quad \text{where } H \text{ pos def } (H = \nabla^2 f(x) \Rightarrow \text{Newton})$$

Plan:

- compute descent direction p_k given x_k
- compute step size $t_k > 0$
- update $x_{k+1} = x_k + t_k p_k$

How to choose t_k ?

- 1) Want sufficient decrease in obj. fun
- 2) Avoid unnecessary small step sizes

1): Simple decrease: $f(x_k + t_k p_k) < f(x_k)$ is not enough:

Example: could imagine method that enforces:

$$f(x_k + t_k p_k) \geq f(x_k) - \frac{1}{2^k} \quad \leq \epsilon$$
$$f(x_{k+1}) \geq f(x_k) - \frac{1}{2^k} \dots \Rightarrow f(x_0) - \sum_{j=0}^k \frac{1}{2^j}$$

steepest method is not good as it decreases f values by at most 2!

Sufficient decrease condition (SDC)

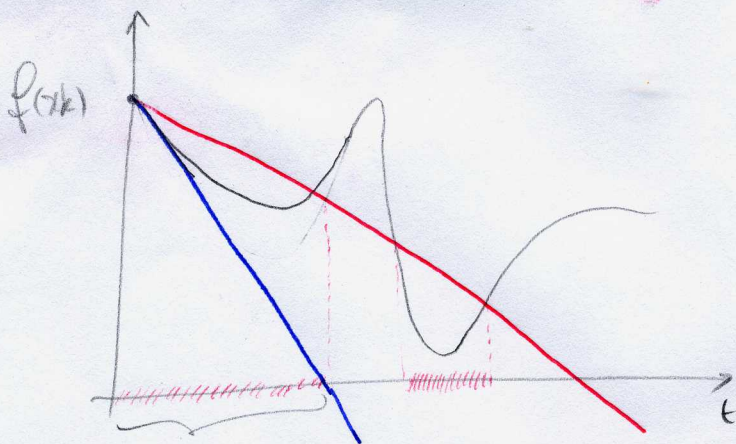
$$f(x_k + t p_k) \leq f(x_k) + c_1 t k \overbrace{\nabla f(x_k)^T p_k}^{< 0}$$

where $0 < c_1 < 1$, $c_1 \sim 10^{-4}$ typical.

motivation: $\varphi(t) = f(x_k + t p_k) \Rightarrow \varphi(0) = f(x_k)$
 $\varphi'(t) = \nabla f(x_k + t p_k)^T p_k$

(SDC) $\Leftrightarrow \varphi(t_k) \leq \varphi(0) + c_1 t_k \varphi'(0)$

\Rightarrow decrease is a fraction of decrease predicted by Taylor's theorem.



maybe problematic that all small values of α are taken!

$$\varphi(0) + \varphi'(0)t$$

$$\varphi(0) + c_1 \varphi'(0)t$$

"less negative" slope

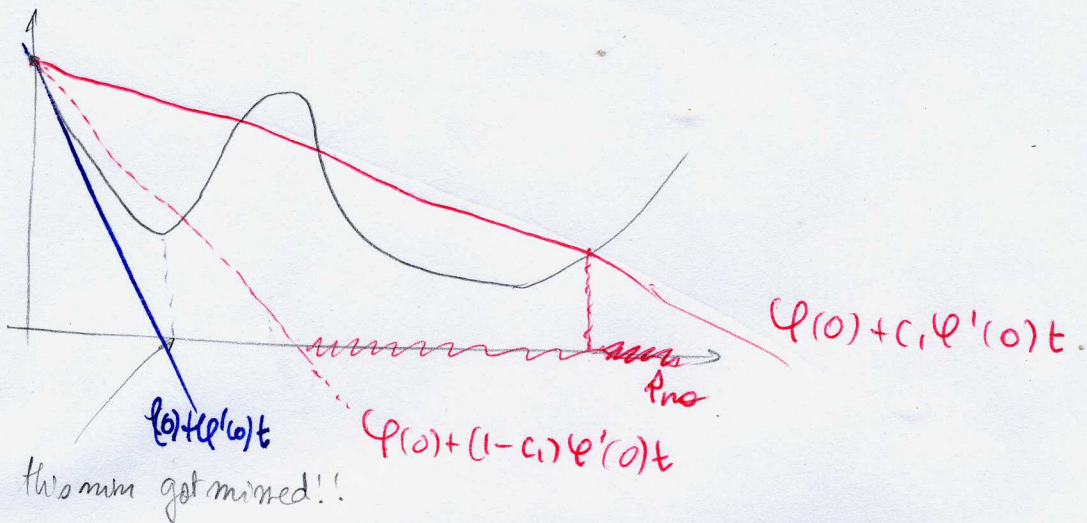
(SDC) is not enough to guarantee convergence: need c_d^0 which ensure t_k is suff large.

4 type of c_d^0 : Wolfe, Strong Wolfe, Goldstein & Backtracking

3) Goldstein cdt^o

$$\begin{cases} f(x_k + t_k p_k) \leq f(x_k) + c_1 t_k \nabla f_k^T p_k \\ f(x_k + t_k p_k) \geq f(x_k) + (1 - c_1) t_k \nabla f_k^T p_k \end{cases}$$

where $c_1 \in (0, \frac{1}{2})$ for two cdt^o to be simultaneously satisfied



Goldstein & Wolfe cdt^o are similar in convergence proof

4) Backtracking (simple)

Let $0 < \beta_1 \leq \beta_2 < 1$

Set $t_k^{(0)}$

for $i = 0 \dots$

if $t_k^{(i)}$ satisfies (SDC) STOP

else select $t_k^{(i+1)} \in [\beta_1 t_k^{(i)}, \beta_2 t_k^{(i)}]$

$\beta_1 = \beta_2 = \frac{1}{2}$ Armijo rule (step size gets halved each time)

Thm If the step size t_k satisfies the Wolfe conditions, Wolfe or Goldstein or Backtracking then:

$\exists c > 0$ (indep of x) s.t.

$$t_k \geq -c \frac{\nabla f_k^T p_k}{\|p_k\|^2} \quad (*)$$

Thm (Convergence of line search methods)

If f is bounded below then:

$$- \sum_{k=0}^{\infty} t_k \nabla f(x_k)^T p_k < \infty \quad (\text{convergent series})$$

$$\Rightarrow \sum_{k=0}^{\infty} \frac{(\nabla f(x_k)^T p_k)^2}{\|\nabla f(x_k)\|^2 \|p_k\|^2} < \infty$$

$\cos^2 \theta_k$; $\theta_k = \text{angle between } \nabla f(x_k) \text{ and } p_k$

In order to get $\nabla f(x_k) \rightarrow 0$ we need to ensure

$$\cos^2 \theta_k \geq \gamma > 0$$

$$\text{If } p_k = -H_k^{-1} \nabla f(x_k)$$

upper bound for num

$$|\nabla f(x_k)^T p_k| = \nabla f(x_k)^T H_k^{-1} \nabla f(x_k) \geq \frac{1}{\lambda_{\max}(H_k)} \|\nabla f(x_k)\|_2^2$$

lower bound for denom

$$\|p_k\|^2 = \|H_k^{-1} \nabla f_k\|^2 = \nabla f_k^T H_k^{-2} \nabla f_k \leq \frac{1}{[\lambda_{\min}(H_k)]^2} \|\nabla f(x_k)\|_2^2$$

$$\Rightarrow \cos^2 \theta_k = \frac{|\nabla f_k^T p_k|}{\|\nabla f_k\|_2 \|p_k\|_2} \geq \left(\frac{\lambda_{\min}(H_k)}{\lambda_{\max}(H_k)} \right)^2 = \text{cond}(H_k)^{-2}$$

need $\text{cond}(H_k)^{-2} \geq \gamma$ for all k

$$\text{cond}(H_k) \leq \sqrt{\gamma^{-1}}$$

Transition to local convergence (p 46)

If:
$$\lim_{k \rightarrow \infty} \frac{\| \nabla^2 f(x_k) p_k + \nabla f(x_k) \|}{\| p_k \|} = 0 \quad (*)$$

- (i) $x_k \rightarrow x^*$ at which 2nd order self optim. cond hold
- (ii) t_k chosen according to rule t_k , Goldstein
- $\Rightarrow \exists \bar{k} \in \mathbb{N}$ s.t. $\forall k > \bar{k} \quad t_k = 1$.

what it means: line search method behaves like Newton method asymptotically, regardless of starting point

\Rightarrow Globally convergent method.

(*) is satisfied if p_k is computed as the solution:

$$(\nabla^2 f(x_k) + \mu_k I) p_k = -\nabla f(x_k) \quad \begin{matrix} \mu_k \geq 0 \\ \mu_k \rightarrow 0 \end{matrix}$$

μ_k guarantees Hessian is pos def but this can be done in various ways by playing with Cholesky factor of Hessian.

Cholesky: $A = LL^T$ (LU factor for symm matrices, very robust factor and reveals whether a matrix is pos def or not).