# 3. Stepsize Selection and Descent

Recall Trust Region steps involve solving the nonconvex minimization (considered last time)

indefinite

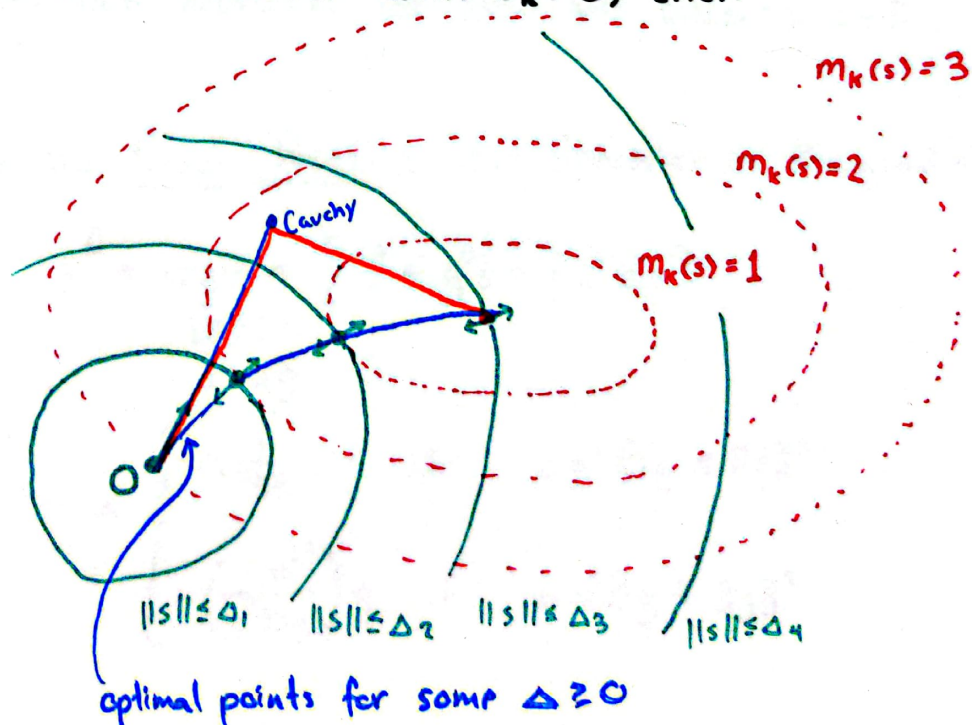$$S_k = \underset{\|s\|_2 \leq \Delta_k}{\text{argmin}} \left\{ \underbrace{f(x_k) + \nabla f(x_k)^T s + \tfrac{1}{2} s^T B_k s}_{m_k(s)} \right\}$$

where $B_k$ models the Hessian
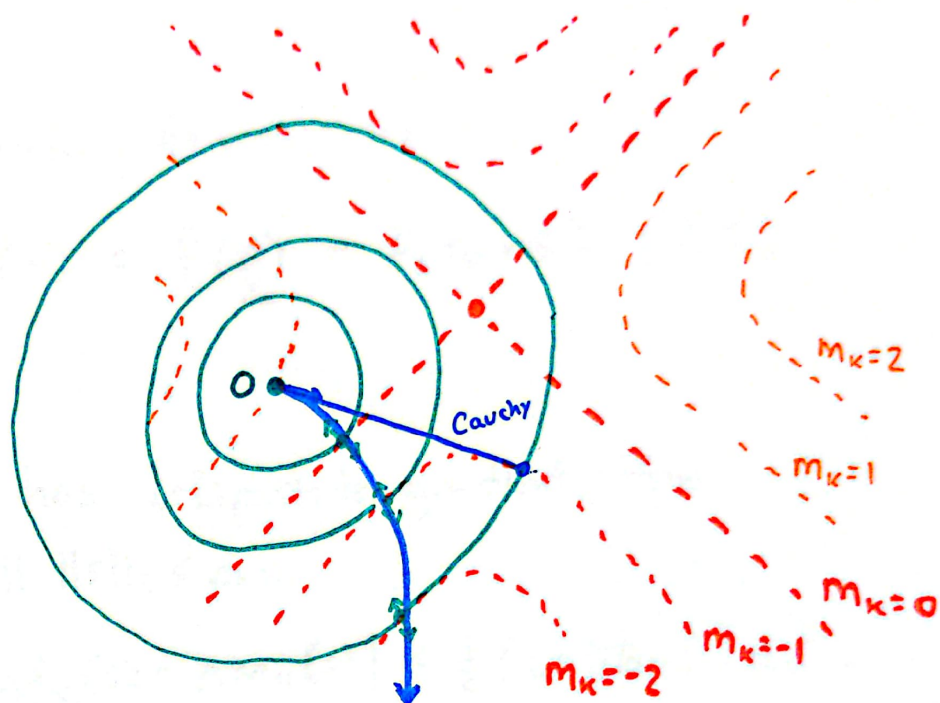and $\Delta_k$ limits the nearby area to search.

## What does the solution look like as we vary $\Delta_k$?

If $m_k(s)$ is convex with $B_k \succ 0$, then



$m_k(s) = 3$

$m_k(s) = 2$

$m_k(s) = 1$

Cauchy

$\|s\| \leq \Delta_1$   $\|s\| \leq \Delta_2$   $\|s\| \leq \Delta_3$   $\|s\| \leq \Delta_4$

optimal points for some $\Delta \geq 0$

If $m_\kappa(s)$ is nonconvex (say $B_\kappa$ indefinite), then

$$m_\kappa(s) = (s_1 - 2)^2 - (s_2 - 1)^2$$



This motivates two heuristics for some (*) decently well...

Define the <u>Cauchy Point</u> as the minimize in the grad direction.

$$s^c = \underset{\substack{\|s\| \leq \Delta \\ s = \alpha g}}{\text{argmin}} \left\{ f + g^\top s + \tfrac{1}{2} s^\top B s \right\}$$

$$= \begin{cases} -\Delta \dfrac{g}{\|g\|} & \text{if} \quad \Delta g^\top B g \leq \|g\|^2 \\[3mm] -\left(\dfrac{\|g\|^2}{g^\top B g}\right) g & \text{if} \quad \Delta g^\top B g \geq \|g\|^2 \end{cases}$$

Define the _Dogleg Path_ as a mixture of gradient and Newton directions (assuming $B \succ 0$):

$$s^{DL}(\tau) = \begin{cases} \tau s^{GD} & \text{if } 0 \le \tau \le 1 \\ \\ s^{GD} + (\tau-1)(s^{N} - s^{GD}) & \text{if } 1 \le \tau \le 2 \end{cases}$$

where $s^{GD} = -\left(\frac{\|g\|^2}{g^T B g}\right) g$ and $s^{N} = -B^{-1}g$.

Pick $s_k$ minimizing $m_k(s^{DL}(\tau))$.


One more heuristic.  Solve in 2D subspace $s^{GD}, s^{N}$,

$$s_k = \underset{\substack{\|s\| \le \Delta \\ s \in \text{span}(-g, -B^{-1}g)}}{\arg\min} f + g^T s + \tfrac{1}{2} s^T B s$$


Back to guaranteeing descent.

Not true that every $\Delta$ gives descent.

Small $\Delta$s work.

Define model objective decrease as

$$\Delta m_K(s) = m_K(0) - \underline{m_K(s)} \quad (>0)$$

and function value decrease as

$$\Delta f_K(s) = f(\underset{x_K}{\bullet}) - f(x_K + s) \quad (\geq 0)$$

**Lemma 1**  If $f$ has $L$-Lipschitz gradient, then
for all $\|s\|_2 \leq \Delta_K$

$$|\Delta f_K(s) - \Delta m_K(s)| \leq \tfrac{1}{2}(L + \|B_K\|)\Delta_K^2.$$

If $f$ has $Q$-Lipschitz Hessian then

$$|\Delta f_K(s) - \Delta m_K(s)| \leq \tfrac{Q}{6}\Delta_K^3 + \frac{\|B_K - \nabla^2 f(x_K)\|}{2}\Delta_K^2$$

Proof.  
$$|\Delta f_K(s) - \Delta m_K(s)| = | f(x_K + s) - (f(x_K) + \nabla f(x_K)^T s + \tfrac{1}{2}s^T B_K s)|$$

$$\leq | f(x_K + s) - (f(x_K) + \nabla f(x_K)^T s)| + \tfrac{1}{2}|s^T B_K s|$$

$$\leq \tfrac{L}{2}\|s\|_2^2 + \frac{\|B_K\|}{2}\|s\|_2^2$$

$$\leq \tfrac{1}{2}(L + \|B_K\|)\Delta_K^2. \qquad \checkmark \quad \overset{\text{Adding and subtracting}}{\swarrow}$$

$$|\Delta f_K(s) - \Delta m_K(s)| \leq | f(x_K + s) - (2^{nd}\text{order model})| + \tfrac{1}{2}|s^T(\nabla^2 f(x_K) - B_K)s|$$

$$\leq \tfrac{Q}{6}\Delta_K^3 + \tfrac{1}{2}\|\nabla^2 f(x_K) - B_K\|\Delta_K^2. \qquad \square$$

## Lemma 2 The Cauchy Point $s^c$ has

$$\Delta m_k(s^c) \geq \tfrac{1}{2}\|\nabla f(x_k)\| \cdot \min\left[\frac{\|\nabla f(x_k)\|}{\|B_k\|}, \Delta_k\right].$$

**Proof.** If $\Delta_k g_k^T B_k g_k \leq \|g_k\|^2$, where $g_k = \nabla f(x_k)$

$$\Delta m_k(s^c) = \Delta\|g_k\| - \tfrac{1}{2}\Delta^2 \frac{g_k^T B_k g_k}{\|g_k\|^2}$$

$$\geq \Delta\|g_k\| - \tfrac{1}{2}\Delta\|g_k\|$$

$$= \tfrac{1}{2}\Delta\|g_k\|.$$

Otherwise $\Delta_k g_k^T B_k g_k > \|g_k\|^2$

$$\Delta m_k(s) = \frac{\|g_k\|^4}{g_k^T B_k g_k} - \tfrac{1}{2}\frac{\|g_k\|^4}{g_k^T B_k g_k}$$

$$= \tfrac{1}{2}\frac{\|g_k\|^4}{g_k^T B_k g_k}$$

$$\geq \tfrac{1}{2}\frac{\|g_k\|^2}{\|B_k\|} \qquad \text{by } g_k^T B_k g_k \leq \|B_k\|\|g_k\|^2.$$

$\square$

Together these give a descent bound

$$|\Delta f_k(s) - \Delta m_k(s)| \leq \tfrac{1}{2}(L + \|B_k\|)\Delta_k^2 \qquad \text{by Lemma 1}$$

$$\Rightarrow \Delta f_k(x_k) \geq \Delta m_k(s) - \tfrac{1}{2}(L + \|B_k\|)\Delta_k^2$$

$$\geq \tfrac{1}{2}\|\nabla f(x_k)\|\min\left\{\frac{\|\nabla f\|}{\|B_k\|}, \Delta_k\right\} - \tfrac{1}{2}(L + \|B_k\|)\Delta_k^2.$$

$$> 0 \quad \text{for small } \Delta > 0.$$

# 4. A Full Trust Region Method

Let's measure how much we trust a step $s$

as $\quad \rho_k(s) := \dfrac{\Delta f_k(s)}{\Delta m_k(s)} \quad$ (Note $\to 1$ as $\Delta_k \to 0$).
$$\Updownarrow$$
$$s \to 0$$

If $\rho_k(s)$ near one or greater, this step is great!

If $\rho_k(s)$ near zero or less, this step is bad!

Picking Thresholds $\quad 0 < \eta_s \le \eta_{vs} \le 1, \quad x_0, \Delta_0,$
$$\overset{\text{``}0.1}{} \qquad \overset{\text{``}0.9}{}$$

we iterate with

for $k = 0, 1, 2, \ldots$

1. Build $m_k(s)$

2. Find $s_k$ minimizing $m_k(s)$ $\quad \begin{pmatrix} \text{at least as well} \\ \text{as Cauchy.} \end{pmatrix}$
   s.t. $\|s\| \le \Delta_k$

3. Compute $\rho_k(s_k)$

4. If $\rho_k(s_k) \ge \eta_{vs}$ $\qquad$ <span style="border:1px solid blue">Very Successful !</span>

   $x_{k+1} = x_k + s_k$
   $\Delta_{k+1} = 2 \Delta_k \quad \leftarrow \gamma_{sv} \cdot \Delta_k$

   Else if $\rho_k(s_k) \ge \eta_s$ $\qquad$ <span style="border:1px solid blue">Success!</span>

   $x_{k+1} = x_k + s_k$
   $\Delta_{k+1} = \Delta_k$

   Else

   $x_{k+1} = x_k$
   $\Delta_{k+1} = \Delta_k / 2.$ $\qquad$ <span style="border:1px solid red">Unsuccessful :c</span>

# 5. Convergence Guarantees

We won't have many unsuccessful steps

since $\rho_k(s_k) = \dfrac{\Delta f_k(s_k)}{\Delta m_k(s_k)}$

$$\geq \frac{\Delta m_k(s_k) - \frac{1}{2}(L + \|B_k\|)\Delta_k^2}{\Delta m_k(s_k)}$$

$$= 1 - \frac{(L + \|B_k\|)\Delta_k^2}{\|\nabla f(x_k)\| \min\left\{ \frac{\|\nabla f\|}{\|B_k\|}, \Delta_k \right\}}$$

$\Rightarrow$ Halving $\Delta_k$ makes $\rho_k$ converge linearly to $1$.

**Claim:** $\Delta_k \to 0$ only if $\|\nabla f(x_k)\| \to 0$.

**Proof.** Suppose all $\|\nabla f(x_k)\| > \varepsilon > 0$. Unsuccessful steps have

$$\eta_s \geq \rho_k = 1 - \frac{(L + \|B_k\|)\Delta_k^2}{\|\nabla f\| \min\left\{ \frac{\|\nabla f\|}{\|B_k\|}, \Delta_k \right\}}$$

$$\Longleftrightarrow \quad 1 - \eta_s \leq \max\left\{ \frac{\beta(L+\beta)\Delta_k^2}{\varepsilon^2}, \frac{(L+\beta)\Delta_k}{\varepsilon} \right\}$$

where $\beta \geq \|B_k\|$.

$$\Longleftrightarrow \quad \Delta_k \geq \min\left\{ \sqrt{\frac{1-\eta_s}{\beta(L+\beta)}}\, \varepsilon, \right.$$

$$\left. \frac{1-\eta_s}{L+\beta}\, \varepsilon \right\}.$$

$\Rightarrow$ At any iteration $\Delta_k \geq \frac{1}{2}\left( \qquad \right). \Rightarrow \Delta_k \not\to 0.$ $\square$

Theorem (Global Convergence, 2018, Curtis, Lubberts, Robinson)
$\quad\quad\quad\quad$ "Concise Complexity Analyses for Trust Region Methods) ← AMS

$$O\left(\tfrac{1}{\varepsilon^2}\right) \text{ rate}, \quad \left(O\left(\tfrac{1}{\varepsilon^{3/2}}\right) \text{ rate with improvements}\right)$$

Theorem $\quad$ If $B_k = \nabla^2 f(x_k)$ $\quad$ (or converging to it),
$\quad\quad\quad$ superlinear convergence

$\quad\quad\quad\quad$ (Nocedal + Wright, Thm 4.9).

# Semester Recap

We have built the machinery and theory for solving

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

for a huge variety of functions $f$.

## Optimality Conditions — What can we locally guarantee and when is this globally meaningful.

## First-Order Optimization — Methods that scale in dimension

- Smooth OPT with optimal acceleration,
- Nonsmooth OPT with subgradients and prox,
- Stochastic/Coordinate Methods with even cheaper per iteration costs,
- Conjugate Gradients and Least Squares.

## Second-Order Methods — Methods that scale in accuracy

- Newton's Method with Quadratic Convergence,
- Quasi Newton (BFGS) with Superlinear Convergence,
- Trust-Regions for Indefinite Local Improvement.

Thank you all for your attention