Recall:   $\nabla^2 f(x_k)\, P_k^N = -\nabla f(x_k) \curvearrowright (*)$

We want:

* robust and efficient in all cases

* $P_k^N$ to be a descent direction:

→ will be true if $\nabla^2 f(x_k)$ is positive definite.

→ else, ...

Consider line-search and trust-region implementations using Newton + CG:

* terminate at negative curvature, or

* modify $\nabla^2 f(x_k)$ to be pos. def.

Also, use Hessian sparsity to get efficient methods.

NB: Some of the most reliable and powerful methods for both small or large are based on Newton methods.

Start with ⊛:

$$\nabla^2 f(x_k)\, P_k^N = -\nabla f(x_k)$$

Then, the residual error is:

$$r_k = \nabla^2 f(x_k)\, P_k + \nabla f(x_k) \sim \text{⊛⊛}$$

for some $P_k \approx P_k^N$.

The problem here is that $cf(\cdot)$ for some constant $c$, will also scale the residual:

$$\underline{\text{New}} \quad r_k' = c\left(\nabla^2 f(x_k)\, P_k + \nabla f(x_k)\right)$$

∴ Eventhough both $f(\cdot)$ and $cf(\cdot)$ have the same stationary points, if we choose to terminate on $|r_k|$, we will stop at different locations for $f(\cdot)$ and $cf(\cdot)$ since $r_k \neq r_k'$.

Instead, we stop for:

$$\|r_k\| \leq \eta_k \|\nabla f(x_k)\|$$

for some sequence, the forcing sequence

$$\eta_1, \eta_2, \dots$$

and $\quad 0 < \eta_k < 1$

We then have convergence according to:

## Thm 6.1 Assume:

* $\nabla f(x)$ is contsly diff'ble in a neighborhood of the minimizer $x^*$.

* $\nabla^2 f(x)$ is positive definite.

Then: $\quad x_{k+1} = x_k + P_k$, and if each

$P_k$ is such that:

$$\|r_k\| \leq \eta_k \|\nabla f(x_k)\|, \quad \eta_k < 1,$$

and if the starting point is sufficiently close to $x^*$, we have convergence.

We will have:

<u>Linear Convergence if:</u>

$$\limsup_{k \to \infty} \frac{\|\nabla f(x_{k+1})\|}{\|\nabla f(x_k)\|} < 1,$$

<u>Quadratic Convergence if:</u>

$$\limsup_{k \to \infty} \frac{\|\nabla f(x_{k+1})\|}{\|\nabla f(x_k)\|^2} = c < \infty, \text{ some } c.$$

<u>Superlinear Convergence if:</u>

$$\limsup_{k \to \infty} \frac{\|\nabla f(x_{k+1})\|}{\|\nabla f(x_k)\|} = 0$$

To get the desired convergence rate, simply set the termination criterion as suggested by theorem 6.2:

<u>Thm 6.2</u>  Assume the setup in thm 6.1. Then:

(I)  $\eta_k \to 0 \implies$ <u>superlinear convergence</u>

(II)  $\eta_k = O(\|\nabla f(x_k)\|) \implies$ <u>quadratic convergence.</u>

The book recommends:

(I) Superlinear convergence:

Use $\eta_k = \min\left(0.5, \sqrt{\|\nabla f(x_k)\|}\right).$

(II) Quadratic Convergence:

Use $\eta_k = \min\left(0.5, \|\nabla f(x_k)\|\right).$

Make sure you can show how (I)+(II) yield the promised convergence...

## Algorithm 6.1 (Line-search Newton-CG)

Given $x_0$

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

$\left\{ \begin{array}{l} \text{Solve } \nabla^2 f(x_k) p = -\nabla f_k \\ \text{using CG, } x^{(0)} = 0. \text{ (see below)} \end{array} \right.$

$\left\{ \begin{array}{l} \text{Compute } \alpha_k \text{ using line-search} \\ \text{that satisfies Wolfe, Goldstein, or} \\ \text{Armijo conditions} \end{array} \right.$

$\qquad x_{k+1} = x_k + \alpha_k p_k$

end

## Modified CG

$x_0 = 0, \quad A = (\nabla^2 f_k), \quad b = -\nabla f_k.$

$r_0 = -b, \quad p_0 = -r_0, \quad k = 0.$

while $\| r_k \| \leq \min\left(0.5, \sqrt{\| \nabla f_k \|}\right) \| \nabla f(x_k) \|$

$\qquad den = p_k^T A p_k$

$\left\{ \begin{array}{l} \text{if } (den \leq 0) \\ \quad \left\{ \begin{array}{l} \text{if } (k > 0) \text{ or } (den \text{ "small"}) \quad \underline{\text{stop}} \\ \underline{\text{else}} \quad \alpha_0 = (r_k^T r_k)/den, \quad x_1 = -\alpha_0 r_0, \quad \underline{\text{stop}} \end{array} \right. \\ \text{end} \end{array} \right.$

$$\alpha_K = r_K^T r_K / den$$

$$x_{K+1} = x_K + \alpha_K p_K$$

$$r_{K+1} = r_K + \alpha_K \boxed{Ap_K}$$

note that this could have been pre-computed with den:

$$\boxed{d} = Ap_K$$
$$den = p_K^T d$$

$$\beta_{K+1} = \frac{r_{K+1}^T r_{K+1}}{r_K^T r_K}$$

$$p_{K+1} = -r_{K+1} + \beta_{K+1} p_K$$

$$K = K+1$$

end.

(* retuns with $p_K$ for the next step *)

---

Algorithm 6.2 (Line search with Preconditioning)

Given $x_0$

for $K = 0, 1, 2, \ldots,$ max_iterations

{Modify $B_K$ to $B_K + E_K$ (discussed next)

{Solve $B_K p_K = -\nabla f(x_K)$ for $p_K$

{Compute $\alpha_K$ using line-search conditions.

$$x_{K+1} = x_K + \alpha_K p_K$$

end

Notes:

· Section 6.3 shows how to compute appropriate $E_k$.

* thm 6.3 says that <u>if</u>:

     ① $\{x \in D : f(x) \leq f(x_0)\}$ <u>compact</u>

     (closed + bdd set that contains its boundary)

     ② $\text{cond}(B_k) = \|B_k\| \, \|B_k^{-1}\| \leq C < \infty$

<u>then</u>

$$\lim_{k \to \infty} \nabla f(x_k) = 0.$$

<u>Proof</u>: Easy!

* $B_k p_k = -\nabla f_k$ can be solved using Gaussian elimination or any standard method since $B_k$ is modified to be positive definite. CG or Gaussian elimination will work.

# 6.3 Hessian Modifications.

We consider:

* $\nabla^2 f(x_k) + \lambda I$, already discussed,
* Controlling selected eigenvalues:

Start from:

$$\nabla^2 f(x_k) = Q \Lambda Q^T \quad \text{as before.}$$

Add $Q \Lambda' Q^T$ to have:

$$Q \Lambda Q^T + Q \Lambda' Q^T = Q \left( \Lambda Q^T + \Lambda' Q^T \right)$$
$$= Q ( \Lambda + \Lambda' ) Q^T$$

The new eigenvalue matrix satisfies:

$$\Lambda + \Lambda' = \text{diag} \left( \lambda_1 + \tau_1, \lambda_2 + \tau_2, \ldots, \lambda_n + \tau_n \right)$$

where:

$$\Lambda = \text{diag} ( \lambda_1, \lambda_2, \ldots, \lambda_n ),$$
$$\Lambda' = \text{diag} ( \tau_1, \tau_2, \ldots, \tau_n ).$$

## Basic idea:

if original eigenvalues are
sufficiently positive:

$$\lambda_i > \varepsilon \geqslant 0 \text{, some } \varepsilon, \text{ then}$$

set $\tau_i = 0$ (don't change them),

## else

### Either:

① Set $\lambda_i + \tau_i = 0 \Rightarrow \lambda_i' = -\tau_i$,

② Flip the sign: $\lambda_i + \tau_i = |\lambda_i|$, or

$$\Rightarrow \lambda_i' = |\lambda_i| + \tau_i$$

③ Make them "sufficiently positive":

$$\text{Set } \lambda_i + \tau_i = \varepsilon \Rightarrow \underline{\underline{\tau_i = \varepsilon - \lambda_i}}.$$

The last option is optimal in the
sense that the change $\Delta A$ to
guarantee that we have

"sufficient positive definiteness" will have $\|\Delta A\|_F$ minimum for ③,

where

$$\|\Delta A\|_F = \sqrt{\sum_{i=1}^{N} \sum_{J=1}^{N} \Delta A_{iJ}^2} \quad . \quad .$$

Using the same notation, for

$$\Delta A = \tau I, \quad \text{set} \quad \tau = \max\left(0, \delta - \lambda_{min}(A)\right)$$

with new matrix $A + \tau I$.

The previous methods: $I, II, III$ are optimal, but they can be very slow.

The textbook describes a fast algorithm in Algorithm 6.5, for computing:

$$PAP^T + E = LDL^T \sim Ⓧ$$

where: $L$ is lower-triangular, $D$ is diagonal, and $E$ is used to make $PAP^T + E$ positive definite.

The book does not explain how to use the factorization in ⊛ to solve $Ax=b$, which is what we need!

We have:

$$PAP^T = LDL^T \implies PA = (LDL^T)P^{-T}$$

Thus: $\qquad Ax = b \implies PAx = Pb$

$$\implies \underbrace{(LDL^T)P^{-T}}x = Pb$$

set as $z_1 = DL^T P^{-T} x$

We have:

Step 1: Solve $Lz_1 = Pb$ for $z_1$.

Step 2: Solve $Dz_2 = z_1$ for $z_2$.

Step 3: Solve $L^T z_3 = z_2$ for $z_3$.

Step 4: Solve $P^{-T} x = z_3$ for $x$.

In other words:
$$z_1 = DL^T P^{-T} x,$$
$$z_2 = L^T P^{-T} x, \text{ and}$$
$$z_3 = P^{-T} x.$$

## Another view:

$$L \underbrace{D \underbrace{L^T \underbrace{P^{-T} x}_{z_3}}_{z_2}}_{z_1} = Pb \qquad \text{and:}$$

$Lz_1 = Pb$ "gives" $z_1$,

$z_1$ "gives" $z_2$ "gives" $z_3$,

and $z_3$ "gives" $x$.

## Computing $P^{-T}$ from $P$, and computing $Pb$.

The book does not show how to compute $P$, $P^T$, or $P^{-T}$!

In algorithm 6.5, we have the step:

"Interchange row and column $J$ and $q$" (of A)

Corresponds to:

$$P_{ij} A' P_{ij}^T$$

where $P_{ij}$ interchanges row $i$ & $J$, where $A'$ is the current value of modified A.

At the output; $A' = P A P^T$.

Now, we can compute $Pb$ by adding

"Interchange rows $J$ and $q$ of $b$"

at this step.

For $P^T$, note that $P_{ij}^T$ can be implemented effectively by simply:

"Interchange column $J$ and $q$!"

Also, note that $P_{ij}^T P_{ij}^T = I$ since interchanging columns $J$ and $q$ twice, brings us back to where we started. $\implies P_{ij}^{-T} = P_{ij}^T$

Also: $P^T = P_{ij(1)}^T P_{ij(2)}^T \cdots P_{ij(n)}^T$

$\underset{\substack{\uparrow \\ \text{1st} \\ \text{iteration}}}{}$ $\qquad\qquad \underset{\substack{\uparrow \\ \text{nth iteration}}}{}$

and

$$P^{-T} = \left(P^T\right)^{-1} = \left( P_{ij(1)}^T P_{ij(2)}^T \cdots P_{ij(n)}^T \right)^{-1}$$

$$= \left( P_{ij(n)}^T \right)^{-1} \cdots \left( P_{ij(1)}^T \right)^{-1}$$

$$= P_{ij(n)}^T \cdots P_{ij(1)}^T$$

Thus, to implement multiplication by $P^{-T}$, simply apply the column exchanges in the reverse order than what was done by algorithm 6.5.

So store the column exchange using:

$$P[J] = q \quad (\text{see below})$$

## Corrected Algorithm 6.5 (see online corrections!)

Given $A, b$

$\gamma = \max\limits_{1 \le i \le n} |a_{ii}| \quad \leftarrow$ largest diagonal element

$\xi = \max\limits_{i \ne j} |a_{ij}| \quad \leftarrow$ largest off-diagonal element

$\delta = u \max(\gamma + \xi, 1) \quad \leftarrow u$ is machine precision

$\boxed{\beta = \max\left(\gamma, \dfrac{\xi}{\sqrt{n^2-1}}, u\right)^{1/2}}$

for $k = 1, 2, \dots, n$

$\quad c_{kk} = a_{kk} \quad \leftarrow$ diagonal

end

for $J = 1, 2, \dots, n$     ($J$th column of $L$).

    Find $q$ so that: $|c_{qq}| \ge |c_{ii}|, \; i = J, J+1, \dots, n$

    Interchange row and column $J$ and $q$ of $A$.
    $\curvearrowleft$ (corrected position of the step!)

    Interchance $J$ and $q$ elements of $b$. $(Pb)$

    $P_J = q \leftarrow$ for inverting: $P^{-T}$.

(Jth column of $L$:)

for $s = 1, 2, \ldots, J-1$

$\quad l_{Js} = c_{Js} / d_s$

end

for $i = J+1, \ldots, n$

$\quad c_{ij} = a_{ij} - \sum_{s=1}^{J-1} l_{Js} c_{is}$

end

$\theta_J = 0$

if $\boxed{J < n}$ ← correct this (see online comments.)

$\quad \theta_J = \max_{J < i \leq n} |c_{iJ}|$

end

$d_J = \max\left\{ |c_{JJ}|, (\theta_J / \beta)^2, \delta \right\}$

if $J < n$

$\quad$ for $i = J+1, \ldots, n$

$\quad\quad c_{ii} = c_{ii} - c_{ij}^2 / d_J$

$\quad$ end

end

end

Note that that the output of the
algorithm keeps:

$$d_J \geq \delta, \quad |m_{iJ}| \leq \beta$$

which keeps eigenvalues away from
zero and $A = MM^T$ not very large.

If you implement the new algorithm,
then:
  ▷ check correctness with example
     6.2.

---

## 6.4 Trust-Region Newton Methods

We always want to compute:

$$\min_{p \in \mathbb{R}^n} m_K(p) \stackrel{def}{=} f_K + \nabla f_K^T p + \tfrac{1}{2} p^T B_K p,$$

such that: $\|p\| \leq \Delta_K$.

For general problems, that are also large,
consider Newton-CG methods.

Use CG to solve: $B_k P_k = -\nabla f_k$, stopping if:

(i) current solution exceeds trust region,

(ii) we have achieved sufficient accuracy (small residuals), or

(iii) we have encountered negative curvature $(P_k^T A f_k < 0)$, (if yes, follow negative curvature to the boundary).

* If we are at a "well-behaved solution", use $\eta_k$ to go to the minimum

* Avoid using the Hessian since the CG method only needs $\nabla^2 f(x_k) v$.

→ Get $\nabla^2 f(x_k) v$ symbolically or

→ Use:
$$\nabla^2 f(x_k) p \approx \frac{\nabla f(x_k + hp) - \nabla f(x_k)}{h}$$

(finite differencing)

<u>Advantages:</u>

* globally convergent:
  - starting with Cauchy point and improving on it.
* No matrix factorizations needed!
* Can be executed in parallel using a matrix times vector routine <u>that</u> executes in parallel.
* Can exploit sparsity in product (avoid multiplying zeros)
* Can move away from nonminimizing points (unlike line-search)

<u>Disadvantages:</u>

* Accepts all directions of negative curvature, even when they lead to an insignificant reduction.

  <u>Eg:</u>  For:  $m(p) = -10^{-3} p_1 - 10^{-4} p_1^2 - p_2^2$
  subject to  $\|p\| \le 1,$   $p = \begin{bmatrix} p_1 \\ p_2 \end{bmatrix}.$

  Correct book!

Now:

$$\nabla m(p) = \begin{bmatrix} -10^{-3} - 2 \times 10^{-4} p_1 \\ -2 p_2 \end{bmatrix}$$

At $p = 0$: $-\nabla m(0) = \begin{bmatrix} 10^{-3} \\ 0 \end{bmatrix}$ steepest descent.

Note:

$$\nabla^2 m(p) = \begin{bmatrix} -2 \times 10^{-4} & 0 \\ 0 & -2 \end{bmatrix}$$

When $p_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$, along steepest descent:

$$p_1^T \begin{bmatrix} -2 \times 10^{-4} & 0 \\ 0 & -2 \end{bmatrix} p_1 = -2 \times 10^{-4} < 0.$$

So, we have negative curvature along $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$. Follow this to boundary of $\|p\| = 1$ to get:

$$m\left(\begin{bmatrix} 1 \\ 0 \end{bmatrix}\right) = -10^{-3} - 10^{-4} \simeq -10^{-3}$$

reduction in model

If we follow $P_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, note:

$$P_2^T \begin{bmatrix} -2 \times 10^{-9} & 0 \\ 0 & -2 \end{bmatrix} P_2 = -2 < 0,$$

we get:

$$m\left(\begin{bmatrix} 0 \\ 1 \end{bmatrix}\right) = -P_2^2 = -1$$

<u>much more reduction</u>.

## Solution I:

Must use "Lanczos method" in references [177], [121].

more robust.

<u>Heuristic</u>: ignore "insignificant" curvatures

## Solution II:

Precondition the Newton-CG method to a better distribution of the eigenvalues:

<u>See</u>: Algorithm 6.6, and get

Code from: MINPACK-2 (NMTR) or LANCELOT.

Final remarks:

For convergence, we set $\eta_k \to 0$ as discussed previously, with Newton-CG, then the <u>trust-region</u> algorithm 4.1 will converge, where the line:

  "Obtain $P_k$ --- "

is replaced by Newton-CG, for positive definite Hessians, $x_k \to x^*$ for these Hessians.

<u>Thm 6.4</u> says that $x_k \to x^*$, for $k \to \infty$ is solved by <u>ignoring the trust region</u>.

∴ Constraint Optimization becomes Unconstraint.

⇒ Trust-regions are for global convergence!!