

- So far: Newton's method
  - local convergence
  - some variants of globalization (backtracking, TR)

- Tonight: Quasi-Newton method

Sometimes, this refers to anything of form

$$x_+ = x - B^{-1} F(x)$$

for some  $B$ . Here assume, successive  $B$ s are obtained by secant updating.

In practice need to globalize.

Here, consider on "local" form above.

- Began with method proposed by W. Davidson (Argonne labs report in 1959) Later suggested by Davidson, Fletcher, Powell (1964)
- Original published in 1<sup>st</sup> edition of SIOPT
- Field gained momentum in 1960's; flourished in 70's matured in early 80's, slackened off after that. How still important but less fundamental.

Recall Newton:  $x_+ = x - J^{-1}(x) F(x)$

Forming  $J(x) \sim \Theta(n^2)$  function evaluation, may be expensive or impossible

so using  $J_S = -F \sim \Theta(n^3)$  arithmetic operations in general

Goal: Develop methods that avoid most of these difficulties.

General methods will require no evaluations of  $J$  and only  $\Theta(n^2)$  arithmetic operations.

Won't enjoy quadratic local convergence but will have superlinear local convergence.

- Recall secant method for  $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$x_+ = x - \left[ \frac{f(x) - f(x_-)}{x - x_-} \right]^{-1} f(x)$$

This has desired form with  $B = \frac{f(x) - f(x_-)}{x - x_-}$

with no  $f'$ -evaluations and superlinear local convergence.

- Suggests that in our general method

$$x_+ = x - B^{-1} F(x)$$

$B_+$  = "next"  $B$  should satisfy

$$B_+ s = y \quad \dots (*) \text{ secant equation}$$

where:  $s = x_+ - x$

$$y = F(x_+) - F(x)$$

for  $n=1$  case, this reduces to secant method

(say  $B_+$  is secant update of  $B$ )

See: This determines  $B_+$  uniquely ( $\Leftrightarrow n=1$ )

What to do when  $n > 1$ ?

- Least-change principle: Make least possible change in  $B$  to get  $B_+$  satisfying  $(*)$
- What does Least-change mean?

Original interpretation - "minimum rank" interpretation

A rank-one update of  $B$  has the form

$$B_+ = B + u v^T \quad \text{for } u, v \in \mathbb{R}^n$$

To satisfy  $(*)$  must have

$$\begin{aligned} y = B_+ s &= B s + u v^T s \\ \Rightarrow v^T s u &= y - B s \Rightarrow u = \frac{y - B s}{v^T s} \end{aligned} \quad \left\{ \Rightarrow B_+ = B + \frac{(y - B s) v^T}{v^T s}$$

work for any  $v$   
such that  $v^T s \neq 0$

Of all choices of  $\vartheta$ , most successful in practice is

$$\vartheta = S$$

$$B_+ = B + \frac{(y - BS)S^T}{S^T S}$$

the Broyden update

introduced in 1965 but still the most successful update for general nonlinear equations.

$$S \leftarrow \lambda S \Rightarrow B_+ = B + \frac{(y - \lambda BS)\lambda S^T}{\lambda^2 S^T S}$$

$$y = F(x + \lambda s) - F(x) = F'(x) \lambda s + O(\lambda)$$

- Shortcomings of minimum-rank interpretation
  - provides no basis for distinguishing Broyden update from others
  - provides no basis for understanding / analyzing method behavior
  - isn't suitable when some "auxiliary conditions" (esp. specific sparsity pattern) are imposed in addition to (\*)
- Modern interpretation: Make the least change in  $B$  as measured in some  $\|\cdot\|$  on  $\mathbb{R}^{n \times n}$  to get  $B_+$  satisfying (\*).  
More precisely

$$B_+ = \underset{\bar{B} s = y}{\operatorname{argmin}} \{ \| \bar{B} - B \| \}$$

denote

$$Q(y, s) = \{ \bar{B} \in \mathbb{R}^{n \times n} : \bar{B} s = y \}$$

and consider

$$\mathcal{H} = \{ H \in \mathbb{R}^{n \times n} : H s = 0 \}$$

$H$ -matrices  
"annihilators" of  $s$

Claim  $\mathcal{H}$  is a subspace of  $\mathbb{R}^{n \times n}$

Why? If  $H_1, H_2$  are in  $\mathcal{H}$  then for scalars  $\alpha_1, \alpha_2$

$$(\alpha_1 H_1 + \alpha_2 H_2) S = \underbrace{\alpha_1 H_1 S}_{=0} + \underbrace{\alpha_2 H_2 S}_{=0} = 0$$

$$\Rightarrow \alpha_1 H_1 + \alpha_2 H_2 \in \mathcal{H}$$

Suppose  $M \in Q(y, s)$  and  $H \in \mathcal{H}$

$$\text{Then } (M+H)S = MS + HS = y + 0$$

$\Rightarrow M+H$  also in  $Q$

If  $M_1, M_2$  are in  $Q$  then

$$(M_1 - M_2)S = M_1 S - M_2 S = y - y = 0$$

$\Rightarrow M_1 - M_2$  is in  $\mathcal{H}$

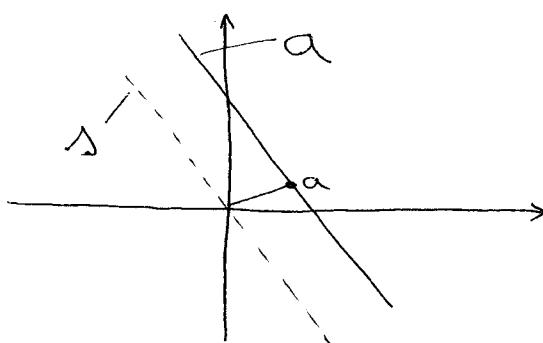
So  $Q(y, s)$  is an affine subspace of  $\mathbb{R}^{n \times n}$  with parallel subspace  $\mathcal{H}$

- Affine subspace - a translate of a subspace.

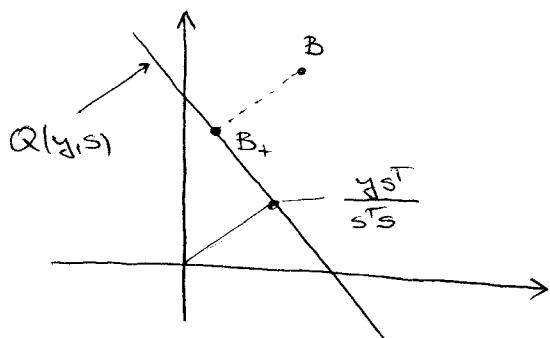
Generate form

$$a = a + \Delta$$

$\uparrow$  vector       $\uparrow$  subspace



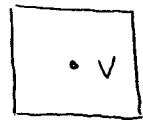
so have



so want  $B_+$  to be the closest element of  $Q$  to  $B$  as measured by  $\|\cdot\|$

- For general  $\|\cdot\|$ , "closest" may not uniquely define  $B_+$
- Example:

$\|\cdot\|_\infty$  on  $\mathbb{R}^n$



all points on boundary  
are equally close to v  
in  $\|\cdot\|_\infty$

$\|\cdot\|_1$  on  $\mathbb{R}^n$



same situation in  $\|\cdot\|_1$ ,

$\Rightarrow$  closest point is not uniquely defined

- For some norms "closest" may uniquely defined but not easy to compute
- Assume  $\|\cdot\|$  is an inner-product norm on  $\mathbb{R}^{n \times n}$   
i.e.  $\|M\| = \sqrt{\langle M, M \rangle}$

Then can characterize  $B_+$  as an orthogonal projection of  $B$  onto  $Q$

- For now, take  $\|\cdot\|$  = Frobenius norm, define  $\mathcal{G}_y$

$$\langle A, B \rangle = \sum_{1 \leq i, j \leq n} A_{ij} B_{ij} = \text{trace}\{A^T B\}$$

so

$$\|A\| = \langle A, A \rangle^{1/2} = \sqrt{\sum_{1 \leq i, j \leq n} A_{ij}^2} = \text{trace}\{A A^T\}$$

$$\text{trace}\{M\} = \sum_i M_{ii} \quad \| \quad$$

Know  $Q(y, s)$  is an affine subspace with parallel subspace  $J\ell$ . Indeed

$$Q(y, s) = \frac{y s^T}{s^T s} + J\ell$$

$$\left[ \begin{array}{l} \text{see } \frac{y s^T}{s^T s} \cdot s = y \end{array} \right]$$

where  $\frac{y^T}{s^T s}$  is the normal element

$$\text{i.e. } \frac{y^T}{s^T s} \in Q(y, s) \cap J\ell^\perp$$

Why? See  $\frac{y^T}{s^T s} \in Q$  and if  $H \in J\ell$  then

$$\left\langle \frac{y^T}{s^T s}, H \right\rangle = \text{trace} \left\langle \frac{y^T}{s^T s} H^T \right\rangle = \text{trace} \left\{ \frac{y^T (\overset{\circ}{H s})^T}{s^T s} \right\} = 0$$

With  $\langle \cdot, \cdot \rangle$  can say  $B_+$  = orthogonal projection of  $B$  onto  $Q$ , i.e. the unique element of  $Q$  such that  $(B_+ - B) \in J\ell^\perp$

$$\text{i.e. } B_+ = P_Q \cdot B$$

where  $P_Q$  = orthogonal projection onto  $Q$

In general, if  $\alpha$  is an affine subspace with normal element  $a_n$  and parabolic subspace  $\Delta$  then

$$P_\alpha \cdot v = a_n + P_\Delta v$$

$$\begin{array}{c} \uparrow \\ \text{orthogonal projection onto } \alpha \end{array} \quad \begin{array}{c} \uparrow \\ \text{orthogonal projection onto } \Delta \end{array}$$

Why?

$$\begin{aligned} (1) \quad P_\alpha (P_\alpha \cdot v) &= a_n + P_\Delta (P_\alpha \cdot v) \\ &= a_n + P_\Delta (a_n + P_\Delta v) \\ &= a_n + \underbrace{P_\Delta a_n}_0 + \underbrace{P_\Delta (P_\Delta v)}_{\text{since } a_n \perp \Delta} \end{aligned}$$

$\parallel P_\Delta v$  since  $\Delta$  is an projection

$$(2) \quad \text{clearly } P_\alpha \cdot v = a_n + P_\Delta v \in \alpha$$

for all  $v \Rightarrow \text{range}(P_\alpha) \leq \alpha$

Also if  $v \in \alpha \Rightarrow v - a_n \in \Delta$

$$\begin{aligned}\Rightarrow P_A(v - av) &= v - av \\ \Rightarrow v &= av + P_A(v - av) \\ v &\in \text{range}(P_A) \\ a &\subseteq \text{range}(P_A) \\ \text{so } \text{range}(P_A) &= a\end{aligned}$$

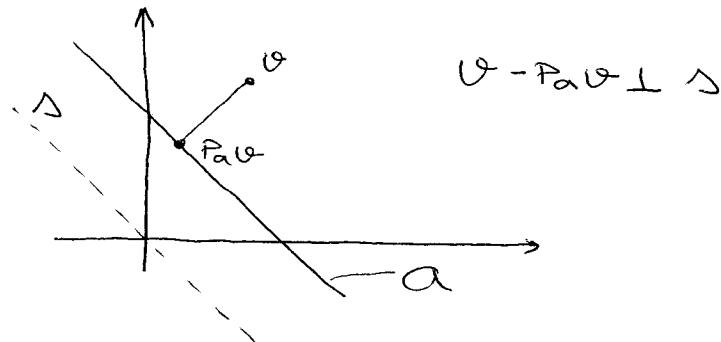
Then (1) & (2) show that  $P_A$  is an algebraic projection  
Also have for any  $v$  and any  $w \in A$

$$\begin{aligned}\langle v - P_A v, w \rangle &= \langle v - av - P_A v, w \rangle \\ &= \underbrace{\langle v - P_A v, w \rangle}_{\parallel 0} - \underbrace{\langle av, w \rangle}_{\parallel 0} = 0\end{aligned}$$

since  $av \perp A$

Since  $P_A$  is orthogonal projection onto  $A$  and therefore  
 $v - P_A v \perp A$

So  $P_A$  is orthogonal projection onto  $\mathbb{Q}$



So have normal element

$$Q(y, s) = \frac{y^T s}{s^T s} + \beta$$

so

$$B_+ = P_Q B = \frac{y^T s}{s^T s} + P_{\mathbb{R}} B$$

In general, if  $A$  is a subspace, then  $P$  is orthogonal projection onto  $A \Leftrightarrow$

- |   |  |
|---|--|
| (1) $P^2 = P$<br>(2) $\text{range } P = A$<br>(3) for all $v$ and all $s \in A$ | $\left. \begin{array}{l} \\ \\ \end{array} \right\} \text{algebraic projection}$<br>$\left. \begin{array}{l} \\ \\ \end{array} \right\} \langle v - P_A v, s \rangle = 0$<br>$\left. \begin{array}{l} \\ \\ \end{array} \right\} \text{orthogonal projection}$ |
|---|--|

Note: Since  $\Delta$  is a subspace

$$(3) \text{ holds } \Leftrightarrow \langle P_S v, w \rangle = \langle v, P_S w \rangle \quad \dots (3')$$

for all  $v, w$

Claim:  $P_{\Delta^{\perp}} B = B \left[ I - \frac{SS^T}{S^T S} \right]$

Why? Verify that (1), (2), (3) or (3') hold

so

$$\begin{aligned} B_+ &= \frac{y S^T}{S^T S} + B \left[ I - \frac{SS^T}{S^T S} \right] \\ &= B + \frac{(y - BS) S^T}{S^T S} \end{aligned}$$

Broyden update again!

Properties of Broyden's method

Doesn't require  $J(x)$  at each step

$$\begin{aligned} x_+ &= x - \bar{B}^{-1} F(x) \quad (\text{But often } B_0 = J(x_0) \text{ in practice}) \\ B_+ &= B + \frac{(y - BS) S^T}{S^T S} \end{aligned}$$

Solving  $BS = -F(x)$  appears to require  $\Theta(n^3)$  arithmetic  
In fact, since  $B_+$  is rank one update of  $B$ , can implement  
only  $O(n^2)$  arithmetic

One way (original way)

Sherman - Morrison - Woodbury formula

If  $A$  nonsingular and  $u, v$  vectors, the  $A + uv^T$   
is nonsingular  $\Leftrightarrow \xi = 1 + v^T A^{-1} u \neq 0$  in which case

$$(A + uv^T)^{-1} = A^{-1} - \frac{1}{\xi} A^{-1} u v^T A^{-1}$$

(Note  $\exists$  generalization to the rank  $k$  update case  
- see Golub - Van Loan )

For Broyden, then can form  $B_0^{-1}$  using  $\Theta(n^3)$  arithmetic

$$\left\{ \begin{array}{l} \text{Factor } B_0 = LU \sim \Theta(n^3) \\ \text{for } j = 1, \dots, n \\ \quad \text{solve } B_0 C_j = e_j \sim \Theta(n^2) \quad \text{for } j^{\text{th}} \text{ column of } I \\ \quad \text{then } B_0^{-1} = (c_1, \dots, c_n) \end{array} \right.$$

Then define

$$B_{k+1}^{-1} = B_k^{-1} - \frac{1}{s_k} B_k^{-1} \frac{(y - B_k s)_k s^T}{s_k^T s_k} B_k^{-1} \quad \text{with } \Theta(n^2) \text{ arithmetic}$$

Preferred way (more complicated but occasionally better numerically)

Form  $B_0 = Q_0 R_0$

Then update  $Q_k$  &  $R_k$  to get  $Q_{k+1}$  &  $R_{k+1}$  in  $\Theta(n^2)$  arith

- Originated by Gill - Murray (1972)
- Can show

Theorem: (Broyden - Dennis - More 1971)

Suppose  $F$  is Lipschitz continuously differentiable at  $x_*$  such that  $F(x_*) = 0$  and  $F'(x_*)$  nonsingular.  
Then for  $x_0$  sufficiently near  $x_*$  and  $B_0$  sufficiently near  $F'(x_*)$ ,  $\{x_k\}$  produced by Broyden's method is well-defined and  $x_k \rightarrow x_*$  superlinearly.

Moreover  $\{B_k\}$  and  $\{B_k^{-1}\}$  are bounded.

- Suppose we want to impose structure\* on  $B_k$ 's reflecting that of  $J(x_k)$

\* e.g. symmetry, particular sparsity etc.

Can extend least-change secant update approach to do this

(1) Suppose have  $\langle \cdot, \cdot \rangle$  and  $\|\cdot\|$  on  $\mathbb{R}^{n \times n}$   
(Frobenius norm or weighted Frobenius norm)

$$\langle A, B \rangle = \text{trace}\{A W B^T\}$$

where  $W = \text{SPD}$  weight matrix

(2) an affine subspace  $\alpha \subseteq \mathbb{R}^{n \times n}$  of matrices having the desired structure

Then, given  $B$  take

$$B_+ = \underset{\bar{B} \in \alpha \cap Q}{\text{argmin}} \|\bar{B} - B\| = \text{least-change secant update at } B \text{ in } \alpha$$

Assuming  $\alpha \cap Q \neq 0$  have

$$B_+ s = y$$

since  $B_+ \in Q$

Also if  $\alpha \cap Q \neq 0$  can show  $\alpha \cap Q$  is an affine subspace: If  $\alpha = \alpha_0 + \Delta$  and  $Q = \frac{y s^T}{s^T s} + \mathcal{J}_\ell$  then

$$\alpha \cap Q = \alpha_0 + \Delta \cap \mathcal{J}_\ell$$

$$\text{so some } n \in (\alpha \cap Q) \cap \mathcal{J}_\ell^\perp$$

Then

$$B_+ = P_{\alpha \cap Q} B = n P_{\alpha \cap Q} B$$

Usually  $\alpha \cap Q \neq 0$

Proposition: If  $J(x) \in \alpha$  for all  $x$  of interest  
then  $\alpha \cap Q \neq 0$

$$\begin{aligned}
 \text{Proof: } & y = F(x_*) - F(x) \\
 &= \int_0^1 \frac{d}{dt} F(x+ts) dt = \underbrace{\int_0^1 J(x+ts) dt \cdot s}_{\in \mathcal{A}}
 \end{aligned}$$

This machinery can actually be applied to derive updates slides 99-102 in short course notes for derivation of the Powell, symmetric Broyden update in which  $\mathcal{A}$  = subspace of symmetric matrices

- See slides 104-107 for many updates that can be derived using this approach.
- See slide 108 for "theorem" giving local superlinear convergence for all of these.
- see slide 109 for recommendations