

Quasi Newton methods:

we have a model function

$$m_k = f_k + \nabla f_k^T p + \frac{1}{2} p^T B_k p$$

B_k is symmetric p.d

B_k could be Hessian, but that isn't
attractive (not p.d? derivatives hard?)

instead, we will update the model

minimizer is: $p_k = -B_k^{-1} \nabla f_k$

next point is:

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

where α_k satisfies Wolfe conds.

(25)

Idea: update B_k

$$M_{k+1} = f_{k+1} + \nabla f_{k+1} P + \frac{1}{2} P' B_{k+1} P$$

→ require gradient of M_{k+1} match
gradient of f_{k+1} at $k, k+1$

$k+1$ already σ_k

k :

$$M_{k+1} \left(-\alpha_k P_k \right) = \nabla f_{k+1} - \alpha_k B_{k+1} P_k = \nabla$$

or

$$\alpha_k B_{k+1} P_k = \nabla f_{k+1} - \nabla f_k$$

write $S_k = \alpha_k P_k$, $y_k = \nabla f_{k+1} - \nabla f_k$

then $B_{k+1} S_k = y_k$

Notice for this to be true $S_k^T y_k > 0$
(But Wolfe cond will ensure this;
which follows from B_{k+1} p.d.)

- 1) No UNIQUE B_{k+1}
 - 2) could require $\|B_{k+1} - B_k\|$ small in ? norm
- careful choice of norm gives

$$B_{k+1} = (I - \gamma_k y_k S_k^T) B_k (I - \gamma_k S_k y_k^T) + \frac{y_k y_k^T}{S_k^T y_k} \cdot \gamma_k$$

where $\gamma_k = \frac{1}{S_k^T y_k}$

This update is DFP.

Notice, we care about $H_k = B_k^{-1}$

$$H_{k+1} = H_k - \frac{H_k y_k y_k^T H_k}{y_k^T H_k y_k} + \frac{S_k S_k^T}{y_k^T S_k}.$$

This is effective, but can do better.

BFGS:

instead of constraining B_k , constrain H_k ($= B_k^{-1}$)

we had $B_k s_k = y_k$

Now, instead, require

$$H_{k+1} y_k = s_k$$

which gives

$$H_{k+1} = (I - \gamma_k s_k y_k^T) H_k (I - \gamma_k y_k s_k^T) + \gamma_k s_k s_k^T$$

where $\gamma_k = \frac{1}{y_k^T s_k}$

This is best, gives superlinear convergence

What is H_0 ?

- options:
- the Hessian inverse
 - Identity
 - Scaled identity
 - $1/\text{diag}$ of Hessian.

Big attraction

limited memory versions

BFGS:

from line search

$$x_{k+1} = x_k - \alpha_k H_k \nabla f_k$$

$$H_{k+1} = V_k^T H_k V_k + \gamma_k S_k S_k^T$$

$$\gamma_k = \frac{1}{y_k^T S_k}$$

$$V_k = (I - \gamma_k y_k S_k^T)$$

Notice we could apply recursively:

$$H_{k+1} = V_k^T \left[(V_{k-1}^T H_{k-1} V_{k-1} + \gamma_{k-1} S_{k-1} S_{k-1}^T) + \gamma_k S_k S_k^T \right] V_k$$

etc.

Notice also:

(30)

- multiply by $S_k S_k^T$ is vector mult.
(ie $S_k S_k^T x = S_k (S_k^T x)$)
- V_k is also easy
 $(I - \gamma_k y_k S_k^T) x = x - \gamma_k (y_k (S_k^T x))$
- recording state is just a question of keeping y_k, S_k, γ_k

IDEA:

- store last m of (y_k, S_k, γ_k)
- throw away oldest
- have an initial H_{k-m} for each step
(Identity; scaled identity; diag)
- This is L-BFGS

METHOD OF CHOICE FOR BIG UNCONSTRAINED PROB