

## Optimization

### Pt. 3 Quasi-Newton Methods

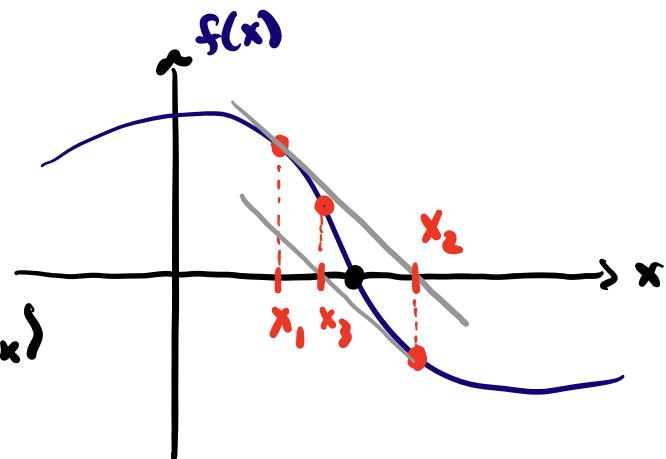
~~Recap~~      Newton's Method (in d=1)

Solve  $f(x) = 0$

Given  $x_0$

for  $k=1, 2, 3, \dots$

$$x_{k+1} = x_k - f(x_k) / f'(x_k)$$



For optimization (critical pts of  $f(x)$ ), run  
Newton to find  $f'(x) = 0$ :

$$x_{k+1} = x_k - f'(x_k) / f''(x_k)$$

$\Rightarrow$  "Second Order" Method

Typically converges quadratically (for  $f(x)=0$ )

(1)  $f''$  continuous

(2)  $f'(x) \neq 0$

(3)  $x_0$  suff. close to root  $x$

## Newton's Method in $d > 1$

$f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ , find  $f(x) = 0$

$$f(x + \Delta x) = f(x) + J_f(x) \Delta x + o(\|\Delta x\|)$$

$$\Rightarrow \Delta x \approx J_f^{-1}(x) [f(x + \Delta x) - f(x)]$$

$\underset{\substack{\Rightarrow \\ = 0 \text{ when} \\ x + \Delta x \text{ is} \\ \text{root of linear approx}}}{}$

Newton

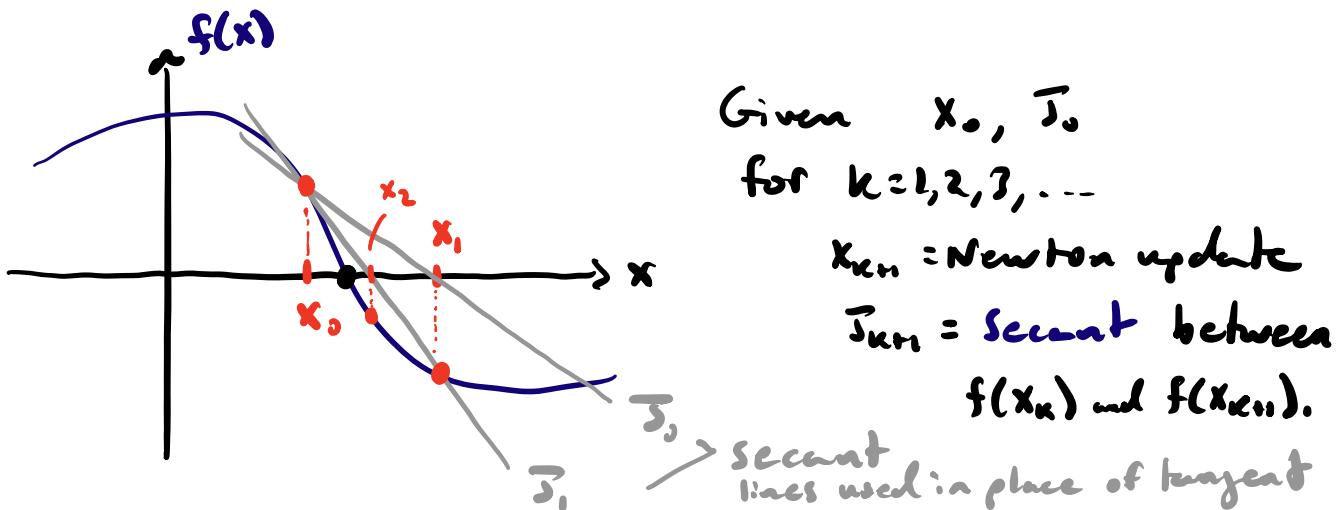
$$\text{Step } \Rightarrow x_{k+1} = x_k + J_f^{-1}(x_k) f(x_k)$$

Note: adapts to non-square systems via  
pseudoinverse or least-squares solve.

## Quasi-Newton Methods

Computing the Jacobian and solving  
 $J_f(x) \Delta x_k = f(x_k)$  at each step is not  
always computationally tractable, e.g.,  
when  $d$  is very large or computing  $f(x)$   
is very expensive and  $J_f$  is not known  
analytically.

Idea: approximate  $\bar{J}_k \approx \bar{J}_S(x_k)$  on the fly, using change updates at each iteration.



Given  $x_0, \bar{J}_0$   
for  $k=1, 2, 3, \dots$

$x_{k+1}$  = Newton update  
 $\bar{J}_{k+1}$  = Secant between  $f(x_k)$  and  $f(x_{k+1})$ .

- If  $x_k \rightarrow x$ , Secant line becomes better and better approx to tangent line
- Secant line requires no explicit knowledge of derivatives of  $f$ .

$$\text{Secant eq. } \underbrace{f(x_{k+1}) - f(x_k)}_{\Delta f_k} = \bar{J}_{k+1} \underbrace{(x_{k+1} - x_k)}_{\Delta x_k}$$

for  $d \geq 1$ , this does not uniquely determine  $\bar{J}_{k+1}$ , so we need to impose more constraints.

"No-change" eq.  $\bar{J}_{k+1} q = \bar{J}_k q$  when  $q^T \Delta x_k = 0$ .

These uniquely determine  $\bar{J}_{k+1}$  from  $\bar{J}_k$ :

$$(\bar{J}_{k+1} - \bar{J}_k)q = u \Delta x^T q \quad \text{for some } u$$

$\Rightarrow$  No. - change satisfied

$$\Delta f_k = [\bar{J}_k + u \Delta x^T] \Delta x = \bar{J}_k \Delta x + u \Delta x^T \Delta x$$

$$u = (\Delta f_k - \bar{J}_k \Delta x) / \Delta x^T \Delta x$$

$\Rightarrow$  Secant satisfied

$$\bar{J}_{k+1} = \bar{J}_k + \underbrace{(\Delta f_k - \bar{J}_k \Delta x)}_{\text{rank-1 update}} \frac{\Delta x^T}{\Delta x^T \Delta x} \rightarrow \bar{J}_k$$

Can show that this choice of  $\bar{J}_k$  satisfies

$$\arg \min \| \bar{J}_{k+1} - \bar{J}_k \|_F$$

s.t. Secant cond.

"Minimal info update" = minimal norm soln to  $\Delta f_k = \bar{J}_{k+1} \Delta x_k$

## Pseudocode (simplest form)

Given  $x_0, \bar{J}_0, f_0 = f(x_0)$

for  $k = 1, 2, 3, \dots$

$$x_{k+1} = x_k - \bar{J}_k^{-1} f_k$$

$$\bar{J}_{k+1} = \bar{J}_k + (\Delta f_k - \bar{J}_k \Delta x_k) \frac{\Delta x_k^T}{\Delta x_k^T \Delta x_k}$$

Since we actually need  $\bar{J}_k^{-1}$  at each step, we can use the **Sherman-Morrison formula** to update  $\bar{J}_{k+1}^{-1}$  from  $\bar{J}_k^{-1}$  fast.

$$\bar{J}_{k+1}^{-1} = \bar{J}_k^{-1} + (\Delta x_k - \bar{J}_k^{-1} \Delta f_k) \frac{\Delta x_k^T \bar{J}_k^{-1}}{\Delta x_k^T \bar{J}_k^{-1} \Delta f_k}$$

This is called **Broyden's first update**.

Broyden's second update comes from applying secant and no-change conditions directly to  $G_{k+1} = \bar{J}_{k+1}^{-1}$ .

$$\text{Secant: } \Delta x_k = G_{k+1} \Delta f_k$$

$$\text{No-change: } G_{k+1} q = G_k q \text{ where } q^T \Delta f_k = 0$$

$$\Rightarrow G_{k+1} = G_k + (\Delta x_k - G_k \Delta f_k) \frac{\Delta f_k^T}{\Delta f_k^T \Delta f_k}$$

The 2<sup>nd</sup> update minimizes

$$\arg \min \|G_{k+1} - G_k\|_F$$

s.t. Secant condition

### BFGS Update

In optimization, consider minimizing

$$f: \mathbb{R}^d \rightarrow \mathbb{R}.$$

Run Newton on  $\nabla f$  using Hessian  $H$ :

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix} \quad H = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_d^2} \end{bmatrix}$$

$\Rightarrow H$  is symmetric when  $f$  is 2x cont. diff.

$\Rightarrow$  At smooth local min,  $H$  is SPD  
(analogous to  $s''(x) > 0$  in 1D)

Notice that Broyden updates are not symmetric. Since we want approximations  $H_K \rightarrow H$  when  $x_K \rightarrow x$  (local min), can we adapt Broyden updates to keep  $H_K$  PSD?

Idea: Use symmetric low-rank update to enforce symmetry while satisfying secant condition.

$$\Delta g_K = \nabla f(x_{K+1}) - \nabla f(x_K), \quad \Delta x_K = x_{K+1} - x_K$$

Secant eq.  $H_{K+1} \Delta x_K = \Delta g_K$

Sym. Rank 2 update:  $H_{K+1} = H_K + \alpha u u^\top + \beta v v^\top$

with  $u = \Delta g_K$  and  $v = H_K \Delta x_K$ , i.e.

$$\Delta H_K = H_{K+1} - H_K = \underbrace{\alpha \Delta g_K (\Delta g_K)^\top}_{\text{symm. type 2 update}} + \underbrace{\beta H_K \Delta x_K (H_K \Delta x_K)^\top}_{\text{symm. type 1 update}}$$

Choose  $\alpha, \beta$  to satisfy secant eqn.

$$\Delta g_K = H_{K+1} \Delta x_K = H_K \Delta x_K + \alpha \Delta g_K (\Delta g_K)^\top \Delta x_K$$

$$+ \beta H_k \Delta x_k \Delta x_k^T H_k \Delta x_k$$

$$\Rightarrow \text{need } \alpha = [(\Delta g_k)^T \Delta x_k]^{-1}, \beta = [\Delta x_k^T H_k \Delta x_k]^{-1}$$

$$\text{to get } \Delta g_k = H_k \cancel{\Delta x_k} + \Delta g_k - H_k \cancel{\Delta x_k} \quad \checkmark$$

$$H_{k+1} = H_k + \frac{\Delta g_k \Delta g_k^T}{\Delta g_k^T \Delta g_k} - \frac{H_k \Delta x_k \Delta x_k^T H_k}{\Delta x_k^T H_k \Delta x_k}$$

We can also use the Sherman-Morrison-Woodbury formula to get a fast rank-2 update formula directly for  $H_k^{-1} \rightarrow H_{k+1}^{-1}$ .

One can also show that BFGS update is PSD under appropriate restrictions on  $\Delta x_k$ .