

CS 450: Numerical Analysis¹

Numerical Optimization

University of Illinois at Urbana-Champaign

¹*These slides have been drafted by Edgar Solomonik as lecture templates and supplementary material for the book “Scientific Computing: An Introductory Survey” by Michael T. Heath ([slides](#)).*

Numerical Optimization

- ▶ Our focus will be on *continuous* rather than *combinatorial* optimization:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to} \quad \mathbf{g}(\mathbf{x}) = \mathbf{0} \quad \text{and} \quad \mathbf{h}(\mathbf{x}) \leq \mathbf{0}$$

- ▶ We consider linear, quadratic, and general nonlinear optimization problems:

Optimality Conditions

- ▶ If \mathbf{x} is an interior point in the feasible domain and is a local minima,

$$\nabla f(\mathbf{x}) = \left[\frac{df}{dx_1}(\mathbf{x}) \quad \cdots \quad \frac{df}{dx_n}(\mathbf{x}) \right]^T = \mathbf{0} :$$

- ▶ *Critical points* \mathbf{x} satisfy $\nabla f(\mathbf{x}) = \mathbf{0}$ and can be minima, maxima, or saddle points:

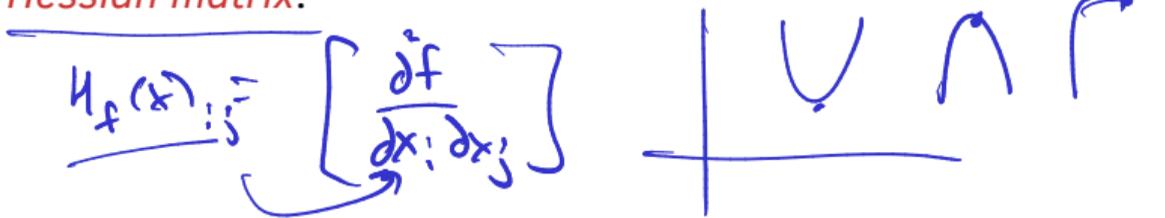
Hessian Matrix

$$\min_x f(x)$$

- ▶ To ascertain whether a critical point x , for which $\nabla f(x) = 0$, is a local minima, consider the *Hessian matrix*:

$$H_f(x) = H_f(x^*)$$

QP



- ▶ If x^* is a minima of f , then $H_f(x^*)$ is positive semi-definite:

sufficient

positive definiteness

$$x^T H x > 0$$

necessary

\Rightarrow local minimum

negative definiteness

\Rightarrow

maximum

indefiniteness

$$\lambda_i(H) > 0$$

$$-\lambda_j(H) < 0$$

\Rightarrow

saddle point

$x+\epsilon$ $f(x+\epsilon) \geq f(x)$
 \Rightarrow local minima

cannot have $\nabla f(x) \neq 0$

$$f(x \pm \epsilon \nabla f(x)) \leq f(x)$$

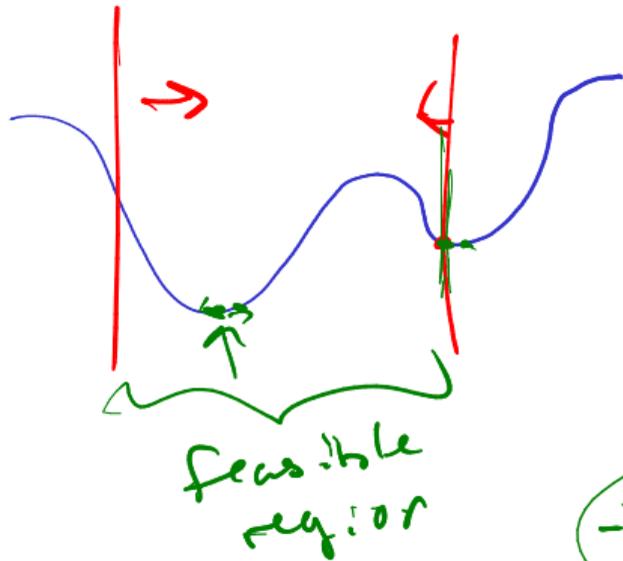
cannot have $s^T H_f(x) s < 0$

$$f(x + \epsilon s) = f(x) + \underbrace{\epsilon \nabla f(x)^T s}_0$$

$$f(x + \epsilon s) < f(x)$$

20

$$+ \frac{\epsilon^2}{2} s^T H_f(x) s + O(\epsilon^3)$$



$$\text{equality } g(x) = 0$$

$$\text{inequality } h(x) \leq 0$$

$$g(x + \epsilon s) = g(x) + \epsilon \underbrace{J_g(x)}_g s + O(\epsilon^2)$$

$$\underbrace{-\nabla f(x)} = \underbrace{J_g^T(x) \lambda}_g \text{ for some } \lambda$$

$$\begin{bmatrix} -3 \\ 0 \\ \vdots \\ 2 \end{bmatrix}$$

$$g(x + \epsilon s) = g(x) + \underbrace{\epsilon J_g(x) s}_g$$

Optimality on Feasible Region Border

- ▶ Given an equality constraint $\mathbf{g}(\mathbf{x}) = \mathbf{0}$, it is no longer necessarily the case that $\nabla f(\mathbf{x}^*) = \mathbf{0}$. Instead, it may be that directions in which the gradient decreases lead to points outside the feasible region:

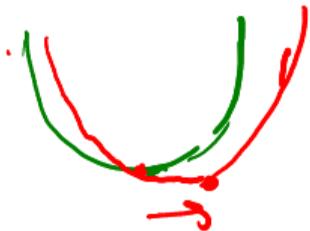
$$\exists \boldsymbol{\lambda} \in \mathbb{R}^n, \quad -\nabla f(\mathbf{x}^*) = \mathbf{J}_g^T(\mathbf{x}^*)\boldsymbol{\lambda}$$


- ▶ Such *constrained minima* are critical points of the Lagrangian function $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})$, so they satisfy:

$$\nabla \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}) = \begin{bmatrix} \nabla f(\mathbf{x}^*) + \mathbf{J}_g^T(\mathbf{x}^*)\boldsymbol{\lambda} \\ \mathbf{g}(\mathbf{x}^*) \end{bmatrix} = \mathbf{0}$$

Sensitivity and Conditioning

- ▶ The condition number of solving a nonlinear equations is $1/f'(x^*)$, however for a minimizer x^* , we have $f'(x^*) = 0$, so conditioning of optimization is inherently bad:

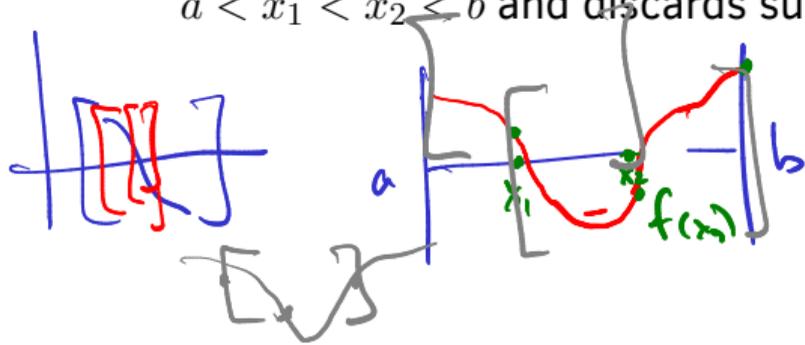


- ▶ To analyze worst case error, consider how far we have to move from a root x^* to perturb the function value by ϵ :

$$\epsilon = f(x^* + h) - f(x^*) = \underbrace{f'(x^*)}_0 h + \frac{1}{2} f''(x^*) h^2 + O(h^3)$$
$$\epsilon \sim h^2 \quad h \sim O(\sqrt{\epsilon})$$

Golden Section Search

- Given bracket $[a, b]$ with a unique local minimum (f is *unimodal* on the interval), *golden section search* considers points $f(x_1), f(x_2)$, $a < x_1 < x_2 < b$ and discards subinterval $[a, x_1]$ or $[x_2, b]$:

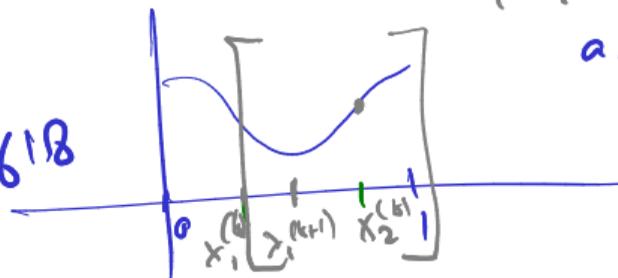


$f(a)$
 $f(x_1)$
 $f(x_2)$
 $f(b)$

which is larger
 $[a, x_2]$ $[b, x_1]$

- Since one point remains in the interval, golden section search selects x_1 and x_2 so one of them can be effectively reused in the next iteration:

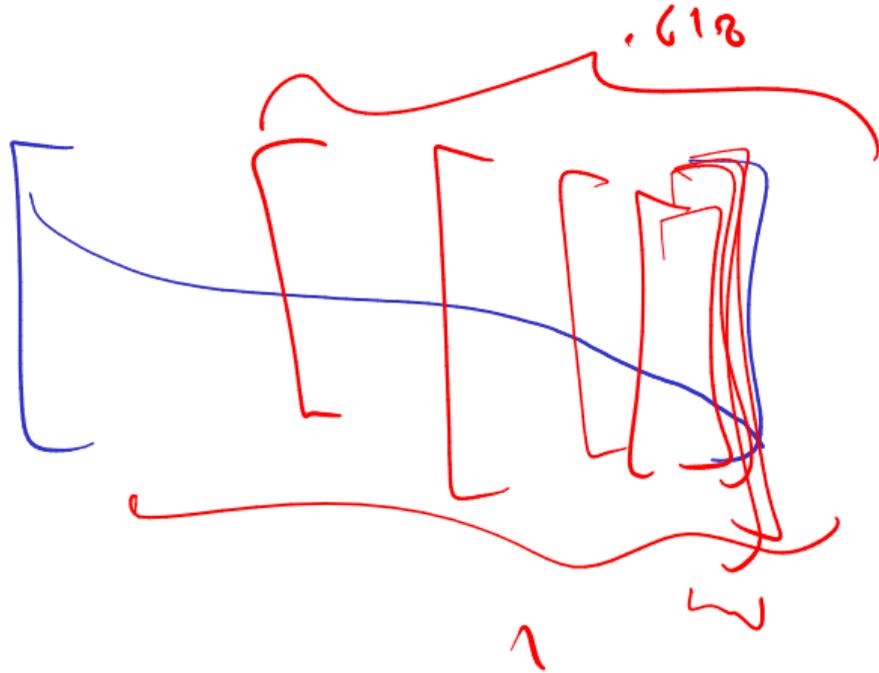
$b - a = 1$
 linear conv rate
 $C = \frac{\sqrt{5} - 1}{2} \approx 0.618$



1 fun-eval/iteration

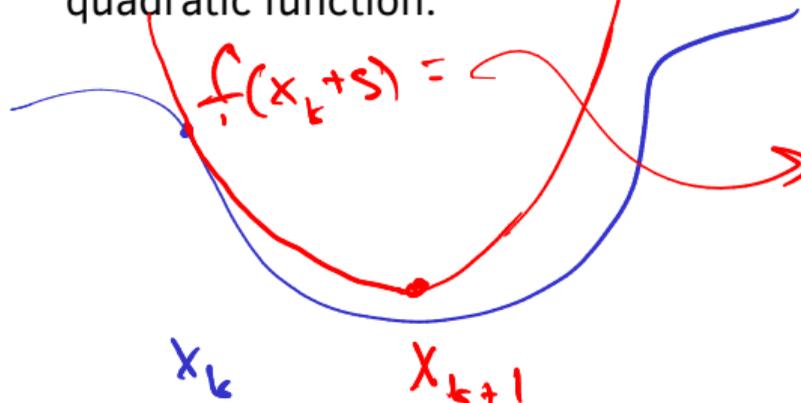
$$a = 0, b = 1$$

$$x_2 = \frac{\sqrt{5} - 1}{2} \quad x_1 = 1 - x_2$$



Newton's Method for Optimization

- At each iteration, approximate function by quadratic and find minimum of quadratic function:



$$\downarrow$$

$$f(x_k) + f'(x_k)s + \frac{1}{2}f''(x_k)s^2$$

$$s = -f'(x_k) / f''(x_k)$$

$$= -h(x) / h'(x)$$

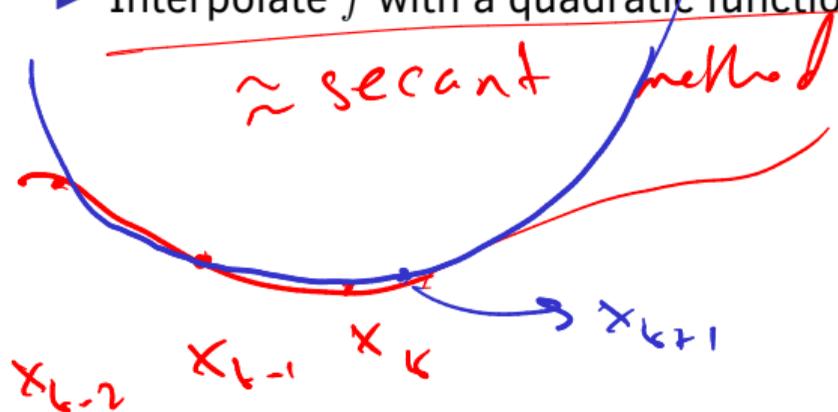
- The new approximate guess will be given by $x_{k+1} - x_k = -f'(x_k) / f''(x_k)$:

= Newton's method

$$\underbrace{f'(x)}_{h(x)} = 0$$

Successive Parabolic Interpolation

- ▶ Interpolate f with a quadratic function at each step and find its minima:

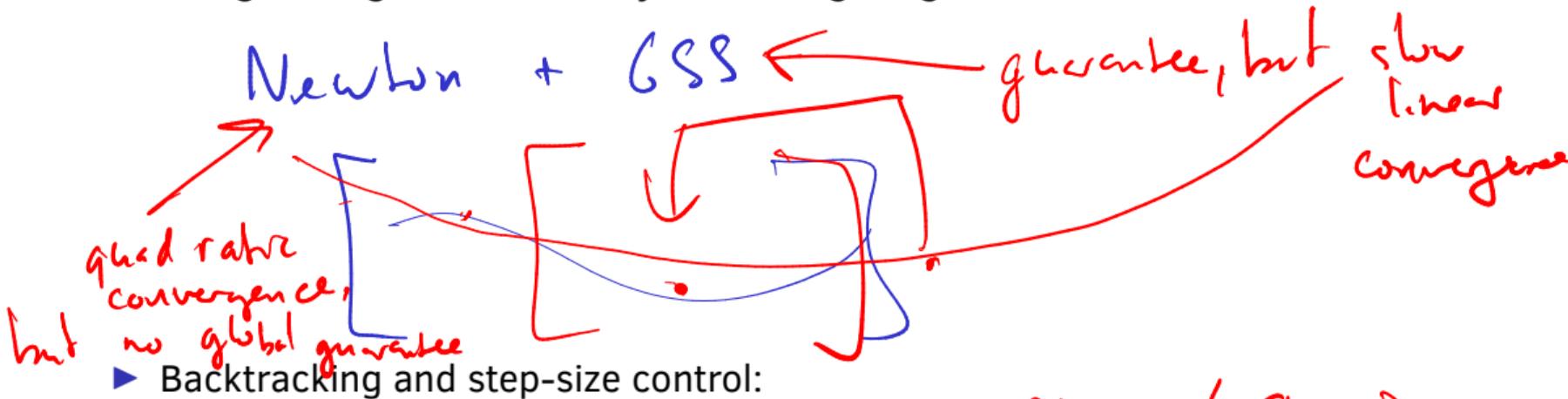


- ▶ The convergence rate of the resulting method is roughly 1.324

superlinear
quadratic

Safeguarded 1D Optimization

- ▶ Safeguarding can be done by bracketing via golden section search:



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

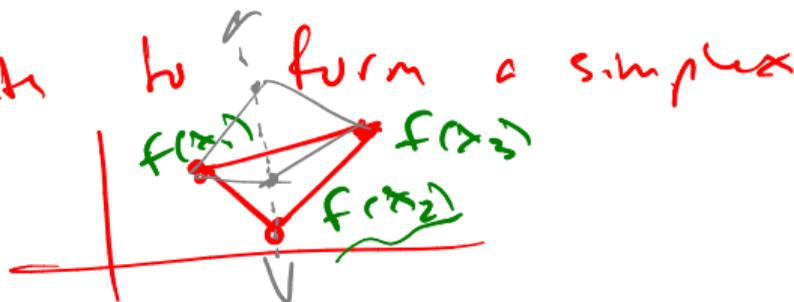
• backtracks α_k if $f(x_{k+1}) > f(x_k)$

General Multidimensional Optimization

- ▶ Direct search methods by simplex (Nelder-Mead):

n-dim opt

n+1 points



- ▶ Steepest descent: find the minimizer in the direction of the negative gradient:

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k)$$

$$\min_{\alpha_k} g(\alpha_k) = f(x_k - \alpha_k \nabla f(x_k))$$

1D optimization

Convergence of Steepest Descent

- ▶ Steepest descent converges linearly with a constant that can be arbitrarily close to 1:

• rate depends on Hessian at x^*

• if gradient changes rapidly, local approximation is poor

- ▶ Given quadratic optimization problem $f(x) = \frac{1}{2}x^T Ax + c^T x$ where A is symmetric positive definite, the error $e_k = x_k - x^*$ satisfies

$$\begin{aligned} \|e_{k+1}\|_A &= e_{k+1}^T A e_{k+1} = \frac{\sigma_{\max} - \sigma_{\min}}{\sigma_{\max} + \sigma_{\min}} \|e_k\|_A \\ &= \frac{\kappa(A) - 1}{\kappa(A) + 1} \|e_k\|_A \end{aligned}$$

Gradient Methods with Extrapolation

- ▶ We can improve the constant in the linear rate of convergence of steepest descent by leveraging extrapolation methods, which consider two previous iterates (maintain momentum in the direction $x_k - x_{k-1}$):

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k) + \beta_k (x_k - x_{k-1})$$

- ▶ The heavy ball method, which uses constant $\alpha_k = \alpha$ and $\beta_k = \beta$, achieves better convergence than steepest descent:

$$\|e_{k+1}\|_A = \frac{\sqrt{5\kappa(A)} - 1}{\sqrt{5\kappa(A)} + 1} \|e_k\|_A$$

Nesterov's gradient optimization

Conjugate Gradient Method

- ▶ The *conjugate gradient method* is capable of making the optimal choice of α_k and β_k at each iteration of an extrapolation method: for $\mathcal{Q} \mathcal{P}$

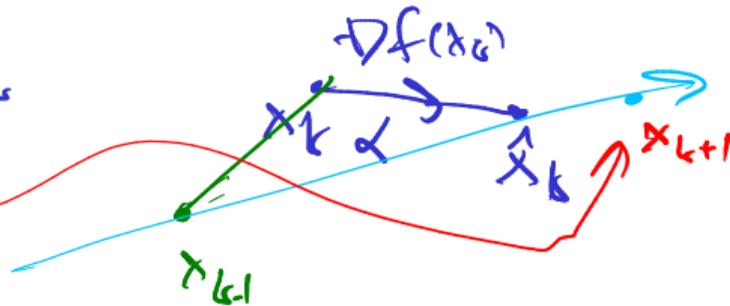
$$\underline{(\alpha_k, \beta_k)} = \underset{\alpha_k, \beta_k}{\operatorname{arg\,min}} f(x_k - \alpha_k \mathcal{P} f(x_k) + \beta_k (x_k - x_{k-1}))$$

Krylov / Lanczos

- ▶ Parallel tangents implementation of the method proceeds as follows

1. steepest descent \hat{x}_k

2. x_{k+1}



Nonlinear Conjugate Gradient

- ▶ Various formulations of conjugate gradient are possible for nonlinear objective functions, which differ in how they compute β below
- ▶ Fletcher-Reeves is among the most common, computes the following at each iteration
 1. Perform 1D minimization for α in $f(\mathbf{x}_k + \alpha \mathbf{s}_k)$
 2. $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{s}_k$
 3. Compute gradient $\mathbf{g}_{k+1} = \nabla f(\mathbf{x}_{k+1})$
 4. Compute $\beta = \mathbf{g}_{k+1}^T \mathbf{g}_{k+1} / (\mathbf{g}_k^T \mathbf{g}_{k+1})$
 5. $\mathbf{s}_{k+1} = -\mathbf{g}_{k+1} + \beta \mathbf{s}_k$

Conjugate Gradient for Quadratic Optimization

$$\nabla f(x) = Ax - b \quad \xrightarrow{\hspace{2cm}} \quad Ax = b$$

- ▶ Conjugate gradient is an optimal iterative method for quadratic optimization,

$$f(x) = \frac{1}{2}x^T Ax - b^T x$$

A is SPD

- ▶ For such problems, it can be expressed in an efficient and succinct form, computing at each iteration

$$r_k = b - Ax_k = -\nabla f(x_k)$$

1. $\alpha = r_k^T r_k / s_k^T A s_k$

2. $x_{k+1} = x_k + \alpha s_k$

3. Compute gradient $r_{k+1} = r_k - \alpha_k A s_k$

4. Compute $\beta = r_{k+1}^T r_{k+1} / (r_k^T r_{k+1})$

5. $s_{k+1} = r_{k+1} + \beta s_k$

matrix vector product

- ▶ Note that for quadratic optimization, the negative gradient $-g$ corresponds to the residual $r = b - Ax$

Krylov Optimization

- ▶ Conjugate Gradient finds the minimizer of $f(x) = \frac{1}{2}x^T Ax - b^T x$ within the Krylov subspace of A :

$$\mathcal{K}(A, b) = \text{span}(b, Ab, \dots, A^{n-1}b) = \text{span}(Q)$$

$$\min_{x \in \mathcal{K}(A, b)} f(x) = \min_{y \in \mathbb{R}^n} = \frac{1}{2} y^T \underbrace{Q^T A Q}_{\text{tridiag}} y - b^T Q y$$

$$x = Qy$$

$$\underline{x = Qy}$$

$$T y^* = \underline{Q^T b} = \|b\|_2 e_1$$

$$y^* = \underline{\|b\|_2^{-1} e_1}$$

Newton's Method

- ▶ Newton's method in n dimensions is given by finding minima of n -dimensional quadratic approximation:

$$\text{= solving } \underline{\nabla f(x) = 0}$$

$$f(x_k + s) \approx \hat{f}(s) = f(x_k) + s^T \nabla f(x_k) + \frac{1}{2} s^T H_f(x_k) s$$

$$0 = \nabla \hat{f}(s) = \nabla f(x_k) + H_f(x_k) s$$

$$H_f(x_k) s = -\nabla f(x_k)$$

$$x_{k+1} = x_k + s = x_k - H_f(x_k)^{-1} \nabla f(x_k)$$

quad. conv.

Quasi-Newton Methods

- ▶ *Quasi-Newton* methods compute approximations to the Hessian at each step:

$$x_{k+1} = x_k - \alpha_k \underline{\underline{B_k^{-1} \nabla f(x_k)}}$$

α_k - line search

- ▶ The *BFGS* method is a secant update method, similar to Broyden's method:

rank-2 update to B_k , using $s_k = x_{k+1} - x_k$, $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$

$$\underline{B_{k+1}} = B_k + \frac{y_k y_k^T}{y_k^T s_k} - \frac{B_k s_k s_k^T B_k}{s_k^T B_k s_k}$$

$O(n^2)$ works / step by Sherman-Morrison

Nonlinear Least Squares

- ▶ An important special case of multidimensional optimization is *nonlinear least squares*, the problem of fitting a nonlinear function $f_{\mathbf{x}}(t)$ so that $f_{\mathbf{x}}(t_i) \approx y_i$:

- ▶ We can cast nonlinear least squares as an optimization problem and solve it by Newton's method:

Gauss-Newton Method

- ▶ The Hessian for nonlinear least squares problems has the form:

- ▶ The *Gauss-Newton* method is Newton iteration with an approximate Hessian:

- ▶ The Levenberg-Marquardt method incorporates Tykhonov regularization into the linear least squares problems within the Gauss-Newton method.

Sequential Quadratic Programming

- ▶ *Sequential quadratic programming* (SQP) corresponds to using Newton's method to solve the equality constrained optimality conditions, by finding critical points of the Lagrangian function $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})$,

- ▶ At each iteration, SQP computes $\begin{bmatrix} \mathbf{x}_{k+1} \\ \boldsymbol{\lambda}_{k+1} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_k \\ \boldsymbol{\lambda}_k \end{bmatrix} + \begin{bmatrix} \mathbf{s}_k \\ \boldsymbol{\delta}_k \end{bmatrix}$ by solving

Barrier Functions

- ▶ *Barrier functions* (*interior point methods*) provide an effective way of working with inequality constraints $\mathbf{h}(x) \leq \mathbf{0}$: