

Math 273a: Optimization  
Newton's methods

Instructor: Wotao Yin  
Department of Mathematics, UCLA  
Fall 2015

some material taken from Chong-Zak, 4th Ed.

## Main features of Newton's method

- Uses both first derivatives (gradients) and second derivatives (Hessian)
- Based on local quadratic approximations to the objective function
- Requires a positive definite Hessian to work
- Converges very quickly near the solution (under conditions)
- Require a lot of work at each iteration:
  - forming the Hessian
  - inverting or factorizing the (approximate) Hessian

## Basic idea

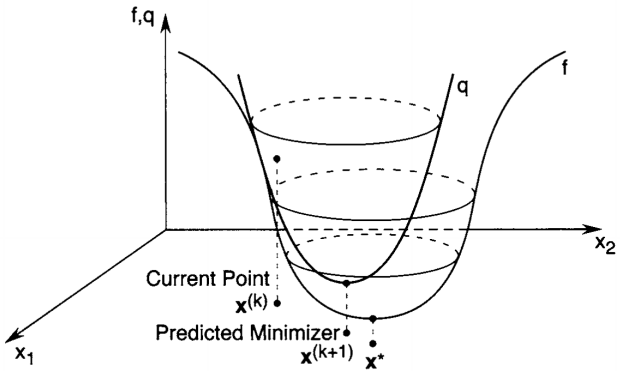
Given the current point  $\boldsymbol{x}^{(k)}$

- construct a quadratic function (known as the quadratic approximation) to the objective function that matches the value and both the first and second derivatives at  $\boldsymbol{x}^{(k)}$
- minimize the quadratic function instead of the original objective function
- set the minimizer as  $\boldsymbol{x}^{(k+1)}$

Note: a new quadratic approximation will be constructed at  $\boldsymbol{x}^{(k+1)}$

Special case: the objective is quadratic

- the approximation is exact and the method returns a solution in one step



## Quadratic approximation

- Assumption: function  $f \in \mathcal{C}^2$ , i.e., twice continuously differentiable
- Apply Taylor's expansion, keep first three terms, drop terms of order  $\geq 3$

$$f(\mathbf{x}) \approx q(\mathbf{x}) := f(\mathbf{x}^{(k)}) + \mathbf{g}^{(k)T}(\mathbf{x} - \mathbf{x}^{(k)}) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(k)})^T \mathbf{F}(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)})$$

where

- $\mathbf{g}^{(k)} := \nabla \mathbf{f}(\mathbf{x}^{(k)})$  is the gradient at  $\mathbf{x}^{(k)}$
- $\mathbf{F}(\mathbf{x}^{(k)}) := \nabla^2 \mathbf{f}(\mathbf{x}^{(k)})$  is the Hessian at  $\mathbf{x}^{(k)}$

## Generating the next point

- Minimizing  $q(\mathbf{x})$  by apply the first-order necessary condition:

$$\mathbf{0} = \nabla q(\mathbf{x}) = \mathbf{g}^{(k)} + \mathbf{F}(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}).$$

- If  $\mathbf{F}(\mathbf{x}^{(k)}) \succ 0$  (positive definite), then  $q$  achieves its unique minimizer at

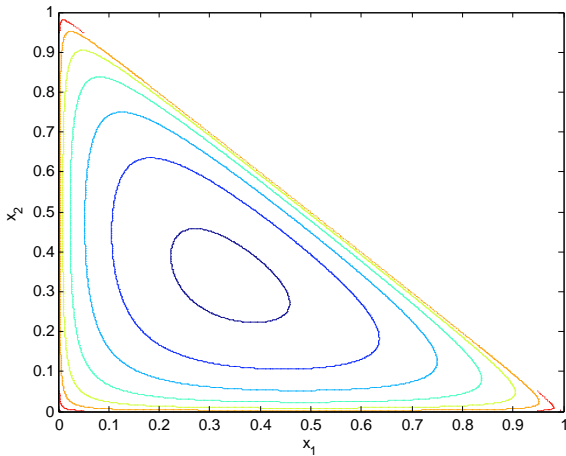
$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}.$$

We have  $\mathbf{0} = \nabla q(\mathbf{x}^{(k+1)})$

- Can be viewed an iteration for solving  $\mathbf{g}(\mathbf{x}) = \mathbf{0}$  using its Jacobian  $\mathbf{F}(\mathbf{x})$ .

## Example

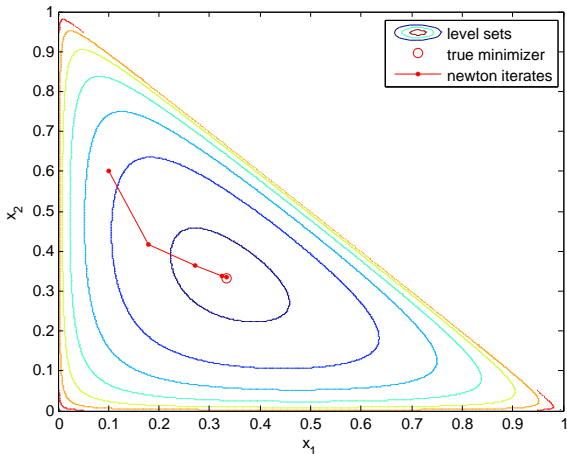
$$f(x_1, x_2) = -\log(1 - x_1 - x_2) - \log(x_1) - \log(x_2)$$



a

## Example

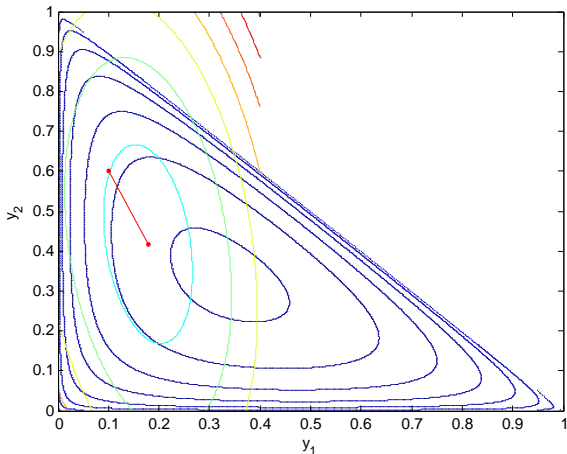
Start Newton's method from  $[\frac{1}{10}; \frac{6}{10}]$





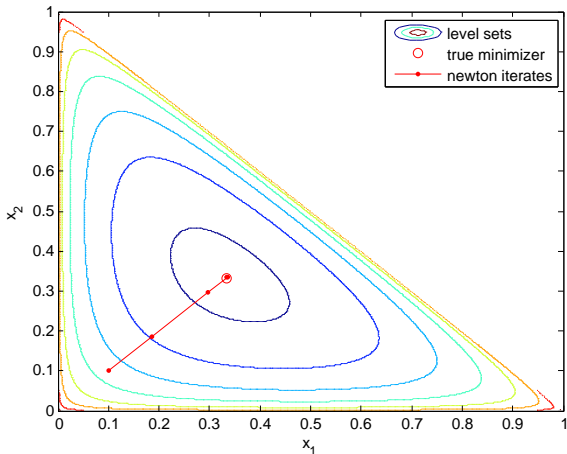
## Example

$f(x_1, x_2)$  and its quadratic approximation  $q(x_1, x_2)$  at  $[\frac{1}{10}; \frac{6}{10}]$  share the same value, gradient and Hessian at  $[\frac{1}{10}; \frac{6}{10}]$ . The new point minimizes  $q(x_1, x_2)$



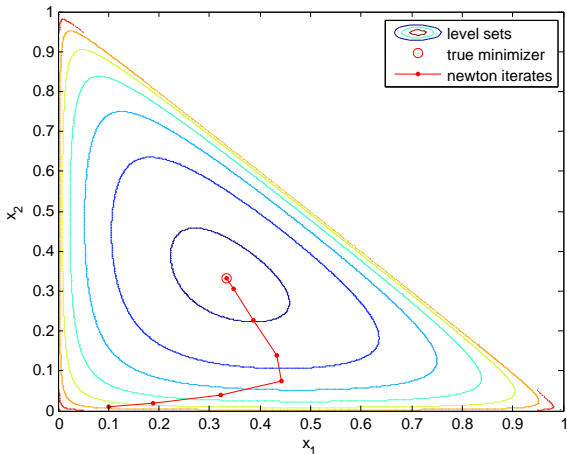
## Example

Start Newton's method from  $[\frac{1}{10}; \frac{1}{10}]$



## Example

Start Newton's method from  $[\frac{1}{10}; \frac{1}{100}]$



## Aware of the drawbacks

Unlike the move along  $-\nabla f(\mathbf{x}^{(k)})$ , where a sufficiently small step size guarantees the objective decrease, Newton's method jumps to a potentially distant point. This makes it vulnerable.

- recall in 1D,  $f'' < 0$  can cause divergence
- even if  $\mathbf{F}(\mathbf{x}^{(k)}) \succ 0$  (positive definite), objective descent is *not* guaranteed

Nonetheless, Newton's method has superior performance when starting near the solution.

## Analysis: quadratic function minimization

- The objective function

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

Assumption:  $\mathbf{Q}$  is symmetric and invertible

$$\mathbf{g}(\mathbf{x}) = \mathbf{Q} \mathbf{x} - \mathbf{b}$$

$$\mathbf{F}(\mathbf{x}) = \mathbf{Q}.$$

- First-order optimality condition  $\mathbf{g}(\mathbf{x}^*) = \mathbf{Q} \mathbf{x}^* - \mathbf{b} = \mathbf{0}$ . So,  $\mathbf{x}^* = \mathbf{Q}^{-1} \mathbf{b}$ .
- Given any initial point  $\mathbf{x}^{(0)}$ , by Newton's method

$$\begin{aligned} \mathbf{x}^{(1)} &= \mathbf{x}^{(0)} - \mathbf{F}(\mathbf{x}^{(0)})^{-1} \mathbf{g}^{(0)} \\ &= \mathbf{x}^{(0)} - \mathbf{Q}^{-1} (\mathbf{Q} \mathbf{x}^{(0)} - \mathbf{b}) \\ &= \mathbf{Q}^{-1} \mathbf{b} \\ &= \mathbf{x}^*. \end{aligned}$$

The solution is obtained in one step.

## Analysis: how fast is Newton's method? (assumption: Lipschitz continuous Hessian near $\mathbf{x}^*$ )

- Let  $\mathbf{e}^{(k)} := \mathbf{x}^{(k)} - \mathbf{x}^*$

### Theorem

Suppose  $f \in \mathcal{C}^2$ ,  $\mathbf{F}$  is Lipschitz continuous near  $\mathbf{x}^*$ , and  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ . If  $\mathbf{x}^{(k)}$  is sufficiently close to  $\mathbf{x}^*$  and  $\mathbf{F}(\mathbf{x}^*) \succ 0$ , then  $\exists C > 0$  such that

$$\|\mathbf{e}^{(j+1)}\| \leq C\|\mathbf{e}^{(j)}\|^2, \quad j = k, k+1, \dots$$

### Just a sketch proof.

Since  $\mathbf{F}$  is Lipschitz around  $\mathbf{x}^*$  and  $\mathbf{F}(\mathbf{x}^*) \succ 0$ , we can have

- $(\mathbf{F}(\mathbf{x}) - c\mathbf{I}) \succ 0$  for some  $c > 0$  for all  $\mathbf{x}$  in a small neighborhood of  $\mathbf{x}^*$ .
- Thus,  $\|\mathbf{F}(\mathbf{x})^{-1}\| < c^{-1}$  for all  $\mathbf{x}$  in the neighborhood.

$\mathbf{e}^{(j+1)} = \mathbf{e}^{(j)} + \mathbf{d}^{(j)}$  where  $\mathbf{d}^{(j)} = -\mathbf{F}(\mathbf{x}^{(j)})^{-1}\mathbf{g}^{(j)}$ .

Taylor expansion near  $\mathbf{x}^{(j)}$ :  $\mathbf{0} = \mathbf{g}(\mathbf{x}^*) = \mathbf{g}(\mathbf{x}^{(j)}) - \mathbf{F}(\mathbf{x}^{(j)})\mathbf{e}^{(j)} + O(\|\mathbf{e}^{(j)}\|^2)$ .

Thus  $\mathbf{e}^{(j+1)} = \mathbf{e}^{(j)} + \mathbf{d}^{(j)} = \mathbf{F}(\mathbf{x}^{(j)})^{-1}O(\|\mathbf{e}^{(j)}\|^2) = O(\|\mathbf{e}^{(j)}\|^2)$ . Argue that  $\mathbf{x}^{(j+1)}, \mathbf{x}^{(j+2)}, \dots$  stay in the neighborhood. □

## Asymptotic rates of convergence

Suppose sequence  $\{\mathbf{x}^k\}$  converges to  $\bar{\mathbf{x}}$ . Perform the ratio test

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{k+1} - \bar{\mathbf{x}}\|}{\|\mathbf{x}^k - \bar{\mathbf{x}}\|} = \mu.$$

- if  $\mu = 1$ , then  $\{\mathbf{x}^k\}$  converges **sublinearly**.
- if  $\mu \in (0, 1)$ , then  $\{\mathbf{x}^k\}$  converges **linearly**;
- if  $\mu = 0$ , then  $\{\mathbf{x}^k\}$  converges **superlinearly**;

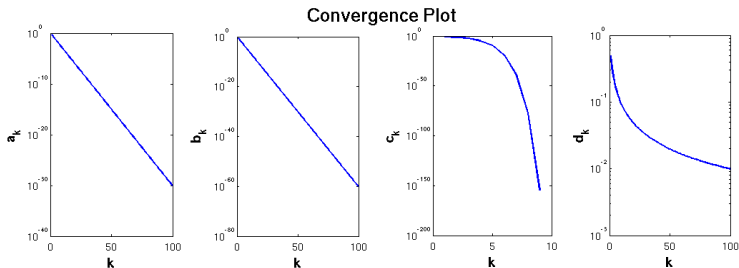
To distinguish superlinear rates of convergence, we check

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{k+1} - \bar{\mathbf{x}}\|}{\|\mathbf{x}^k - \bar{\mathbf{x}}\|^q} = \mu > 0$$

- if  $q = 2$ , it is **quadratic convergence**;
- if  $q = 3$ , it is **cubic convergence**;
- $q$  can be non-integer, e.g., 1.618 for the secant method ...

## Example

- $a_k = 1/2^k$
- $b_k = 1/4^{\lfloor k/2 \rfloor}$
- $c_k = 1/2^{2^k}$
- $d_k = 1/(k + 1)$



“semilogy” plots (wikipedia)



## Another example

Let  $C = 100$ .

$k$	$(1/k^2)$	$Ce^{-k}$	$Ce^{-k^{1.618}}$	$Ce^{-k^2}$
1	1.0e0	3.7e2	3.7e2	3.7e2
3	3.3e-1	5.0e1	2.7e0	1.2e-1
5	2.0e-1	6.7e-0	1.3e-3	1.4e-8
7	1.4e-1	9.1e-1	7.6e-8	5.2e-19
9	1.1e-1	1.2e-1	6.4e-13	6.6e-33

Comments:

- the constant  $C$  is not important in superlinear convergence
- even with a big  $C$ , higher-order convergence will quickly catch up
- the constant  $C$  is more important in lower-order convergence
- superlinear convergence is shockingly fast!

## Analysis: descent direction

### Theorem

If the Hessian  $\mathbf{F}(\mathbf{x}^{(k)}) \succ 0$  (positive definite) and  $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$ , then the search direction

$$\mathbf{d}^{(k)} = -\mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}$$

is a descent direction, that is, there exists  $\bar{\alpha} > 0$  such that

$$f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) < f(\mathbf{x}^{(k)}), \quad \forall \alpha \in (0, \bar{\alpha}).$$

### Proof.

Let  $\phi(\alpha) := f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$ . Then  $\phi'(\alpha) := \nabla f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})^T \mathbf{d}^{(k)}$ . Since  $\mathbf{F}(\mathbf{x}^{(k)}) \succ 0$  and  $\mathbf{g}^{(k)} \neq \mathbf{0}$ , we have  $\mathbf{F}(\mathbf{x}^{(k)})^{-1} \succ 0$

$$\phi'(0) := \nabla f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})^T \mathbf{d}^{(k)} = -\mathbf{g}^{(k)T} \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)} < 0.$$

Finally, apply first-order Taylor expansion to  $\phi(\alpha)$  to get the result. □

## Two more issues with Newton's method

### Hessian evaluation:

- When the dimension  $n$  is large, obtain  $F(\mathbf{x}^{(k)})$  can be computationally expensive
- We will study quasi-Newton methods to alleviate this difficulty (in a future lecture)

### Indefinite Hessian:

- When the Hessian is not positive definite, the direction is not necessarily descending.
- There are simple modifications.

## Modified Newton's method

### Strategy:

- use  $F(x^{(k)})$  if  $F(x^{(k)}) \succ \mathbf{0}$  and  $\lambda_{\min}(F(x^{(k)})) > \epsilon$ ; *otherwise*,
- use  $\hat{F}(x^{(k)}) = F(x^{(k)}) + E$  so that  $\hat{F}(x^{(k)}) \succ \mathbf{0}$  and  $\lambda_{\min}(\hat{F}(x^{(k)})) > \epsilon$ .

**Method 1 (Greenstadt):** replace any tiny or negative eigenvalues by

$$\delta = \max\{\epsilon_{\text{machine}}, \epsilon_{\text{machine}} \|H\|_{\infty}\}$$

where  $\|H\|_{\infty} = \max_{i=1, \dots, n} \sum_{j=1}^n |h_{ij}|$ . This is computationally expensive.

**Method 2 (Levenberg-Marquardt):** Was proposed for least-squares but works here. Replace

$$\hat{F} \leftarrow F + \gamma I.$$

It shifts *every eigenvalue* of  $F$  up by  $\gamma$ .

## Modified Newton's method

### Method 3 (advanced topic: modified Cholesky / Gill-Murray):

Any symmetric matrix  $A \succ 0$  can be factored as

$$A = \bar{L}\bar{L}^T \quad \text{or} \quad A = LDL^T,$$

where  $L$  and  $\bar{L}$  are lower triangular,  $D$  is positive diagonal, and  $L$  has ones on its main diagonal.

### Properties of the Cholesky factorization:

- Very useful in solving linear systems of equations. Reduces a system to two backsolves.
- If  $A \not\succeq 0$  (indefinite but still symmetric),  $D$  has zero or negative element(s) on its diagonal.

## Modified Newton's method

- The factorization is **stable** if  $A$  is positive definite. (Small errors in  $A$  will not cause large errors in  $L$  or  $D$ .) Example:

$$\begin{bmatrix} 1 & a \\ a & 1 \end{bmatrix} = \begin{bmatrix} 1 & \\ a & 1 \end{bmatrix} \begin{bmatrix} 1 & \\ & 1 - a^2 \end{bmatrix} \begin{bmatrix} 1 & a \\ & 1 \end{bmatrix}.$$

- If  $A \leftarrow A + \mathbf{v}\mathbf{v}^T$ , the factorization can be updated to a product form (avoiding the factorization from scratch, which is more expensive).
- If  $A$  is sparse, Cholesky with pivots keeps  $L$  sparse with moderately more zeros
- The cost is  $n^3/6 + O(n^2)$ , roughly half of Gaussian elimination.

## Modified Newton's method

**Forsgren, Gill, Murray:** perform *pivoted* Cholesky factorization. That is, permute the matrix at each step to pull the largest remaining diagonal element to the pivot position

**The effect:** postpone the modification and keeps it as small as possible.

When no acceptable element remains

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{11} & \\ \mathbf{L}_{21} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{D}_1 & \\ & \mathbf{D}_2 \end{bmatrix} \begin{bmatrix} \mathbf{L}_{11}^T & \mathbf{L}_{21}^T \\ & \mathbf{I} \end{bmatrix},$$

replace  $\mathbf{D}_2$  (not necessarily diagonal!) by a *positive definition* matrix and complete the factorization.

No extra work if the Cholesky factorization is taken in the outer-product form.

The Cholesky factorization also tells if the current point is a minimizer or a saddle point.

## Modified Newton's method for saddle point

Suppose we are at a saddle point  $\bar{x}$ . Then  $\bar{g} = \nabla f(\bar{x}) = \mathbf{0}$  and 2nd-order approximation

$$f(\bar{x} + \mathbf{d}) \approx q(\mathbf{d}) := f(\bar{x}) + \underbrace{\bar{\mathbf{g}}^T \mathbf{d}}_{=0} + \frac{1}{2} \mathbf{d}^T \mathbf{F}(\bar{x}) \mathbf{d}.$$

How do we descend?

**Greenstadt:** pick  $\mathbf{d} = \sum_{i: \lambda_i < 0} \alpha_i \mathbf{u}_i$ , where  $\alpha_i > 0$  and  $(\lambda_i, \mathbf{u}_i)$  are eigen-pairs of  $\mathbf{F}(x)$ .  $\mathbf{d}$  is a positive linear combination of the *negative curvature* directions.

Then,  $\mathbf{d}^T \mathbf{F}(\bar{x}) \mathbf{d} < 0$  and  $q(\mathbf{d}) < f(\bar{x})$ .

**Cholesky method:** recall  $\mathbf{D}_2$  correspond to the negative curvatures. If entry  $d_{ij}$  of  $\mathbf{D}_2$  has the largest absolute value among all entries of  $\mathbf{D}_2$ , pick

$$\mathbf{L}^T \mathbf{d} = \mathbf{e}_i - \text{sign}(d_{ij}) \mathbf{e}_j.$$

Then,  $\mathbf{d}^T \mathbf{F}(x) \mathbf{d} = \mathbf{d}^T \mathbf{L} \mathbf{D} \mathbf{L}^T \mathbf{d} < 0$ .



## Overview of the Gauss-Newton method

- A modification to Newton's method, solves nonlinear least squares, very popular
- Pros: second derivatives are no longer computed
- Cons: does not apply to general problems
- Can be improved by line search, Levenberg-Marquardt, etc.

## Nonlinear least squares

- Given functions  $r_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $i = 1, \dots, m$
- The goal is to find  $\mathbf{x}^*$  so that  $r_i(\mathbf{x}) = 0$  or  $r_i(\mathbf{x}) \approx 0$  for all  $i$ .
- Consider the nonlinear least-squares problem

$$\underset{\mathbf{x}}{\text{minimize}} \quad \frac{1}{2} \sum_{i=1}^m (r_i(\mathbf{x}))^2 .$$

- Define  $\mathbf{r} = [r_1, \dots, r_m]^T$ . Then we have

$$\underset{\mathbf{x}}{\text{minimize}} f(\mathbf{x}) = \frac{1}{2} \mathbf{r}(\mathbf{x})^T \mathbf{r}(\mathbf{x}).$$

- The gradient  $\nabla f(\mathbf{x})$  is formed by components

$$(\nabla f(\mathbf{x}))_j = \frac{\partial f}{\partial x_j}(\mathbf{x}) = \sum_{i=1}^m r_i(\mathbf{x}) \frac{\partial r_i}{\partial x_j}(\mathbf{x})$$

- Define the Jacobian of  $\mathbf{r}$

$$\mathbf{J}(\mathbf{x}) = \begin{bmatrix} \frac{\partial r_1}{\partial x_1}(\mathbf{x}) & \cdots & \frac{\partial r_1}{\partial x_n}(\mathbf{x}) \\ \vdots & \cdots & \vdots \\ \frac{\partial r_m}{\partial x_1}(\mathbf{x}) & \cdots & \frac{\partial r_m}{\partial x_n}(\mathbf{x}) \end{bmatrix}$$

Then, we have

$$\nabla f(\mathbf{x}) = \mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x})$$

- The Hessian  $F(\mathbf{x})$  is symmetric matrix. Its  $(k, j)$ th component is

$$\begin{aligned}\frac{\partial^2 f}{\partial x_k \partial x_j} &= \frac{\partial}{\partial x_k} \left( \sum_{i=1}^m r_i(\mathbf{x}) \frac{\partial r_i}{\partial x_j}(\mathbf{x}) \right) \\ &= \sum_{i=1}^m \left( \frac{\partial r_i}{\partial x_k}(\mathbf{x}) \frac{\partial r_i}{\partial x_j}(\mathbf{x}) + r_i(\mathbf{x}) \frac{\partial^2 r_i}{\partial x_k \partial x_j}(\mathbf{x}) \right)\end{aligned}$$

- Let  $S(\mathbf{x})$  be formed by  $(k, j)$ th components

$$\sum_{i=1}^m r_i(\mathbf{x}) \frac{\partial^2 r_i}{\partial x_k \partial x_j}(\mathbf{x})$$

- Then, we have  $F(\mathbf{x}) = \mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}) + S(\mathbf{x})$
- Therefore, Newton's method has the iteration

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \underbrace{(\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}) + S(\mathbf{x}))^{-1}}_{\mathbf{F}(\mathbf{x})^{-1}} \underbrace{\mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x})}_{\nabla f(\mathbf{x})}$$

# The Gauss-Newton method

- When the matrix  $\mathbf{S}(\mathbf{x})$  is ignored in some applications to save computation, we arrive at the Gauss-Newton method

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \underbrace{(\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}))^{-1}}_{(\mathbf{F}(\mathbf{x}) - \mathbf{S}(\mathbf{x}))^{-1}} \underbrace{\mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x})}_{\nabla f(\mathbf{x})}$$

- A potential problem is that  $\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}) \not\approx \mathbf{0}$  and  $f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)})$ .

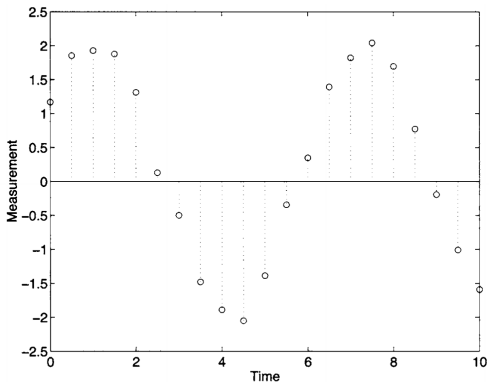
Fixes: line search, Levenberg-Marquardt, and Cholesky/Gill-Murray.

## Example: nonlinear data-fitting

- Given a sinusoid

$$y = A \sin(\omega t + \phi)$$

- Determine parameters  $A$ ,  $\omega$ , and  $\phi$  so that the sinusoid best fits the observed points:  $(t_i, y_i)$ ,  $i = 1, \dots, 21$ .



- Let  $\mathbf{x} := [A, \omega, \phi]^T$  and

$$r_i(\mathbf{x}) := y_i - A \sin(\omega t_i + \phi)$$

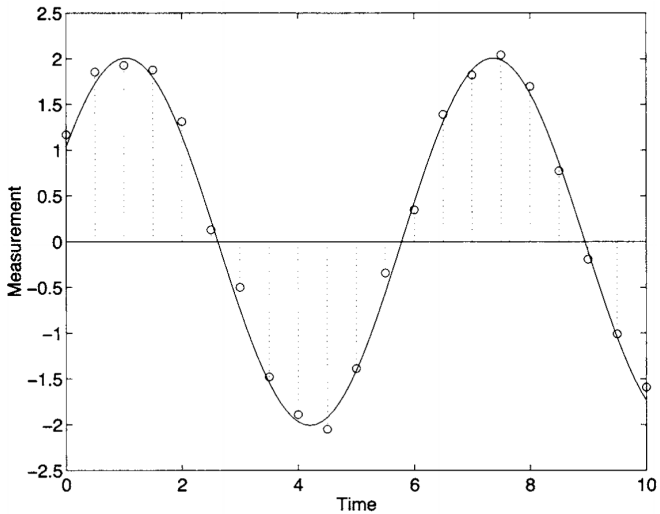
- Problem

$$\text{minimize } \sum_{i=1}^{21} \underbrace{(y_i - A \sin(\omega t_i + \phi))}_{r_i(\mathbf{x})}^2$$

- Derive  $\mathbf{J}(\mathbf{x}) \in \mathbb{R}^{21 \times 3}$  and apply the Gauss-Newton iteration

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - (\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}))^{-1} \mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x})$$

- Results:  $A = 2.01$ ,  $\omega = 0.992$ ,  $\phi = 0.541$ .





## Conclusions

Although Newton's method has many issues, such as

- the direction can be ascending if  $\mathbf{F}'(\mathbf{x}^{(k)}) \neq 0$
- may not ensure descent in general
- must start close to the solution,

Newton's method has the following strong properties:

- one-step solution for quadratic objective with an invertible  $Q$
- second-order convergence rate near the solution if  $\mathbf{F}'$  is Lipschitz
- a number of modifications that address the issues.