

CS 450: Numerical Analysis¹

Nonlinear Equations

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](#)).*

Solving Nonlinear Equations

- ▶ Solving (systems of) nonlinear equations corresponds to root finding:
 - ▶ $f(x^*) = 0$ *univariate nonlinear function*
 - ▶ $f(\mathbf{x}^*) = 0$ *multivariate, scalar-valued nonlinear function*
 - ▶ $\mathbf{f}(\mathbf{x}^*) = \mathbf{0}$ *multivariate, vector-valued nonlinear function*
 - ▶ Algorithms for root-finding make it possible to solve systems of nonlinear equations and employ a similar methodology to finding minima in optimization.
- ▶ Main algorithmic approach: find successive roots of local linear approximations of \mathbf{f} :

- ▶ *Newton's method for univariate functions starting at point x_k finds root of $h(\delta x) = f(x_k) + f'(x_k)\delta x \approx f(x_k + \delta x)$, so*

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

- ▶ *Newton's method for multivariate functions starting at point \mathbf{x}_k finds root of $\mathbf{h}(\delta \mathbf{x}) = \mathbf{f}(\mathbf{x}_k) + \mathbf{J}_f(\mathbf{x}_k)\delta \mathbf{x} \approx \mathbf{f}(\mathbf{x}_k + \delta \mathbf{x})$, with Jacobian $(\mathbf{J}_f(\mathbf{x}))_{ij} = \frac{\delta f_i}{\delta x_j}(\mathbf{x})$, so*

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \delta \mathbf{x} = \mathbf{x}_k - \mathbf{J}_f^{-1}(\mathbf{x}_k)\mathbf{f}(\mathbf{x}_k)$$

Nonexistence and Nonuniqueness of Solutions

- ▶ Solutions do not generally exist and are not generally unique, even in the univariate case:

Consider functions that are strictly greater than zero or have many zeros.

- ▶ Solutions in the multivariate case correspond to intersections of hypersurfaces:

*The zeros of each equation define a **hypersurface** in \mathbb{R}^n , in the linear case, these are **hyperplanes**. Intersections of hypersurfaces for many equations, define the solutions, which are roots of all equations.*

*Consider that two curves can intersect at many points in space. Two hypersurfaces in three-dimensional space may not intersect or may have multiple **curves of intersection**.*

Conditions for Existence of Solution

- ▶ *Intermediate value theorem* for univariate problems:

*If for $a < b$, $\text{sign}(f(a)) \neq \text{sign}(f(b))$ and f is continuous, then $[a, b]$ is a **bracket** that contains a root,*

$$\exists x^* \in [a, b], \quad f(x^*) = 0.$$

- ▶ A function has a unique *fixed point* $g(x^*) = x^*$ in a given closed domain if it is *contractive* and contained in that domain,

$$\|g(x) - g(z)\| \leq \gamma \|x - z\|$$

- ▶ *Contained implies that in the domain S , for any $x \in S$, $g(x) \in S$, while contractive implies that the function is Lipschitz continuous in S .*
- ▶ *When solving for a root of f , can define various fixed point functions g , so that their solution $g(x^*) = x^*$ provides a root of f , $f(x^*) = 0$, the simplest being $g(x) = f(x) + x$.*

Conditioning of Nonlinear Equations

- ▶ Generally, we take interest in the absolute rather than relative conditioning of solving $f(x) = 0$:
 - ▶ *The sensitivity of solving a nonlinear equation is the ratio of magnitudes of the perturbation to the root and perturbation to the function values.*
 - ▶ *It makes sense to consider absolute perturbations to f , since any nonzero perturbation to function values is infinite in magnitude relative to $f(x^*) = 0$.*
- ▶ The *absolute condition number* of finding a root x^* of f is $1/|f'(x^*)|$ and for a root x^* of f it is $\|\mathbf{J}_f^{-1}(x^*)\|$:
 - ▶ *If we change the value of f by at most δf at any point in the function while maintaining continuity, the root will shift by at most $|\delta f|/|f'(x^*)|$ assuming $|\delta f|$ is sufficiently small.*
 - ▶ *This relationship is the converse of conditioning in function evaluation, where a perturbation to input x , results in a perturbation of at most $\kappa_{abs}(f) = |f'(x)|$ larger to the function value.*

Multiple Roots and Degeneracy

- ▶ If x^* is a root of f with **multiplicity** m , its $m - 1$ derivatives are also zero at x^* ,

$$f(x^*) = f'(x^*) = f''(x^*) = \dots = f^{(m-1)}(x^*) = 0.$$

Proof: for some function $t^{(0)}(x)$, we have that

$$f(x) = (x - x^*)^m t^{(0)}(x),$$

$$f'(x) = m(x - x^*)^{m-1} t^{(0)}(x) + (x - x^*)^m t^{(0)'}(x)$$

$$\equiv (x - x^*)^{m-1} t^{(1)}(x),$$

$$f^{(m-1)}(x) = (x - x^*) t^{(m-1)}(x),$$

where $t^{(i)} = (m - i + 1)t^{(i-1)}(x) - (x - x^)t^{(i-1)'}(x)$.*

- ▶ Increased multiplicity affects conditioning and convergence:
 - ▶ When a root x^* not **simple**, i.e. $m > 1$, then $f'(x^*) = 0$, so the problem of finding that root is ill-posed as $1/|f'(x^*)| = \infty$.
 - ▶ In practice, this means we have multiple roots at the same x^* which are impossible to distinguish and may need to reformulate problem/algorithms.

Bisection Algorithm

- ▶ Assume we know the desired root exists in a bracket $[a, b]$ and $\text{sign}(f(a)) \neq \text{sign}(f(b))$:
 - ▶ *if f is continuous, by the intermediate value theorem, the bracket contains a root*
 - ▶ *can proceed to narrow interval to find a root*
 - ▶ *one caveat is that multiple roots may exist in $[a, b]$*
 - ▶ *another caveat is that we've imposed a restrictive condition, it can be difficult to find two points where the function has opposite sign*
- ▶ Bisection subdivides the interval by a factor of two at each step by considering $f(c_k)$ at $c_k = (a_k + b_k)/2$:

$$[a_{k+1}, b_{k+1}] = \begin{cases} [c_k, b_k] & : \text{sign}(f(a_k)) = \text{sign}(f(c_k)) \\ [a_k, c_k] & : \text{sign}(f(b_k)) = \text{sign}(f(c_k)) \end{cases}$$

Rates of Convergence

- ▶ Let x_k be the k th iterate and $e_k = x_k - x^*$ be the error, bisection obtains *linear convergence*, $\lim_{k \rightarrow \infty} \|e_k\| / \|e_{k-1}\| \leq C$ where $C < 1$:

In bisection, working with the natural error bound given by bracket size,

$$e_k = b_k - a_k = \frac{1}{2}(b_{k-1} - a_{k-1}) = \frac{1}{2}e_{k-1},$$

so bisection achieves linear convergence with $C = 1/2$. With linear convergence, error $e_k \leq \epsilon$ is achieved after $O(\log_C(1/\epsilon))$ steps.

- ▶ r th order convergence implies that $\|e_k\| / \|e_{k-1}\|^r \leq C$
 - ▶ Convergence of order $r > 1$ (*superlinear*) implies that the number of digits of correctness increases by a factor of r at each step.
 - ▶ For $r > 1$, error $e_k \leq \epsilon$ is achieved after $O(\log_r(\log(1/\epsilon)))$ steps.
 - ▶ Having achieved superlinear convergence ($r > 1$), methods differ only by constant factors in complexity.

Convergence of Fixed Point Iteration

- ▶ Fixed point iteration: $x_{k+1} = g(x_k)$ is locally linearly convergent if for $x^* = g(x^*)$, we have $|g'(x^*)| < 1$:

Note that,

$$e_{k+1} = x_{k+1} - x^* = g(x_k) - g(x^*).$$

By taking the Taylor expansion of g at x^ to represent $g(x_k)$, we can observe*

$$\begin{aligned} e_{k+1} &= g'(x^*)(x_k - x^*) + O((x_k - x^*)^2) \\ &= g'(x^*)e_k + O((x_k - x^*)^2), \end{aligned}$$

where the asymptotic term decreases as x_k approaches x^ .*

- ▶ It is quadratically convergent if $g'(x^*) = 0$:

Taking the same Taylor expansion, the leading term is now zero and we obtain

$$\begin{aligned} e_{k+1} &= g''(x^*)(x_k - x^*)^2/2 + O((x_k - x^*)^3) \\ &= g''(x^*)e_k^2/2 + O((x_k - x^*)^3). \end{aligned}$$

- ▶ Newton's method is derived from a *Taylor series* expansion of f at x_k :

$$f(x) = \underbrace{f(x_k) + f'(x_k)(x - x_k)}_{\text{secant line approximation}} + (1/2!)f''(x_k)(x - x_k)^2 + \dots$$

- ▶ Newton's method is *quadratically convergent* if started sufficiently close to x^* so long as $f'(x^*) \neq 0$:
 - ▶ Newton's method corresponds to the fixed point function $g(x) = x - f(x)/f'(x)$.
 - ▶ It achieves quadratic convergence since $g'(x^*) = 0$,

$$g'(x^*) = 1 - \underbrace{f'(x^*)/f'(x^*)}_1 + \underbrace{f(x^*)f''(x^*)/f'(x^*)^2}_0.$$

Secant Method

- ▶ The *Secant method* approximates $f'(x_k) \approx \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}$.
 - ▶ Usually, this method is the cheapest approximation possible, since function values $f(x_k)$ and $f(x_{k-1})$ are already available.
 - ▶ Approximation quality depends on magnitude $f(x_k) - f(x_{k-1})$ and $x_k - x_{k-1}$.
 - ▶ If the two points are far apart, the derivative approximation may be bad locally, while if they are very close we have to take care in handling cancellation.
 - ▶ A well-chosen finite-difference step at each x_k provides a more robust approximation, but requires another function evaluation.
- ▶ The convergence of the Secant method is *superlinear* but not quadratic:
 - ▶ The error will now depend on the previous two errors, since we are using the previous two points.
 - ▶ In simplified form, the error at the k th iteration satisfies $e_k \leq e_{k-1}e_{k-2}$.
 - ▶ Note that $\log(e_k) = \log(e_{k-1}) + \log(e_{k-2})$ is the Fibonacci sequence, which grows at a rate of $r = (1 + \sqrt{5})/2$.
 - ▶ Thus the (negative) exponent of the error increases by roughly a factor of r at each step, i.e. the convergence rate is r .

Nonlinear Tangential Interpolants

- ▶ Secant method uses a linear interpolant based on points $f(x_k), f(x_{k-1})$, could use more points and higher-order interpolant:

Have points $(x_0, f(x_0)), \dots, (x_k, f(x_k))$ can fit polynomial to $p(x_i) = f(x_i)$ for some subset of points $x_i \in S \subseteq \{x_0, \dots, x_k\}$.

- ▶ Quadratic interpolation (Muller's method) achieves convergence rate $r \approx 1.84$:

Quadratic interpolation requires three points x_{k-2}, x_{k-1} , and x_k .

- ▶ Inverse quadratic interpolation resolves the problem of nonexistence/nonuniqueness of roots of polynomial interpolants:
 - ▶ *Interpolate quadratic polynomial q so that $q(f(x_i)) = x_i$ for $i \in \{k, k-1, k-2\}$.*
 - ▶ *Approximate root simply computed by evaluating interpolant at zero*
 $x_{k+1} = q(0)$.

Achieving Global Convergence

- ▶ Hybrid bisection/Newton methods:

*Given a bracket (interval), can proceed with bisection until bracket is small then switch to Newton. Alternatively, can attempt Newton, check if it stays within bracket (*safeguard*) and proceed with change only if it does.*

- ▶ Bounded (damped) step-size:

Newton's method gives us a direction. Decreasing the step size in that direction trades off convergence rate for reliability. We will study how step sizes can be chosen in more detail in the context of optimization.

Systems of Nonlinear Equations

- ▶ Given $\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}) \ \cdots \ f_m(\mathbf{x})]^T$ for $\mathbf{x} \in \mathbb{R}^n$, seek \mathbf{x}^* so that $\mathbf{f}(\mathbf{x}^*) = \mathbf{0}$
 - ▶ \mathbf{x}^* must simultaneously set to zero all components of \mathbf{f} :
 $f_1(\mathbf{x}^*) = \cdots = f_m(\mathbf{x}^*) = 0$.
 - ▶ We focus on the case of $m = n$, so that the number of equations matches the number of unknowns.
- ▶ At a particular point \mathbf{x} , the *Jacobian* of \mathbf{f} , describes how \mathbf{f} changes in a given direction of change in \mathbf{x} ,

$$\mathbf{J}_{\mathbf{f}}(\mathbf{x}) = \begin{bmatrix} \frac{df_1}{dx_1}(\mathbf{x}) & \cdots & \frac{df_1}{dx_n}(\mathbf{x}) \\ \vdots & & \vdots \\ \frac{df_m}{dx_1}(\mathbf{x}) & \cdots & \frac{df_m}{dx_n}(\mathbf{x}) \end{bmatrix}$$

Our local linear approximation is given by

$$\mathbf{f}(\mathbf{x} + \delta\mathbf{x}) \approx \mathbf{f}(\mathbf{x}) + \mathbf{J}_{\mathbf{f}}(\mathbf{x})\delta\mathbf{x},$$

note that when $m = 1$ the Jacobian corresponds to the gradient of f .

Multivariate Newton Iteration

- Fixed-point iteration $\mathbf{x}_{k+1} = \mathbf{g}(\mathbf{x}_k)$ achieves local convergence so long as $|\lambda_{\max}(\mathbf{J}_g(\mathbf{x}^*))| < 1$ and quadratic convergence if $\mathbf{J}_g = \mathbf{O}$:

Taking a multivariate Taylor expansion of $\mathbf{g}(\mathbf{x}_k)$ at center \mathbf{x}^ we get*

$$\begin{aligned}
 \mathbf{e}_k &= \mathbf{x}^* - \mathbf{x}_{k+1} \\
 &= \mathbf{g}(\mathbf{x}^*) - \mathbf{g}(\mathbf{x}_k) \\
 &= \underbrace{\mathbf{J}_g(\mathbf{x}^*)}_{\mathbf{O}}(\mathbf{x}^* - \mathbf{x}_k) \\
 &\quad + \frac{1}{2} \left[\mathbf{H}_g^{(1)}(\mathbf{x}^*) \cdot (\mathbf{x}^* - \mathbf{x}_k) \quad \dots \quad \mathbf{H}_g^{(n)}(\mathbf{x}^*) \cdot (\mathbf{x}^* - \mathbf{x}_k) \right] (\mathbf{x}^* - \mathbf{x}_k) \\
 \|\mathbf{e}_{k+1}\| &= O(\max_i \|\mathbf{H}_g^{(i)}(\mathbf{x}^*)\| \cdot \underbrace{\|\mathbf{x}^* - \mathbf{x}_k\|}_{\mathbf{e}_k}^2) \\
 &= O(\|\mathbf{e}_k\|^2)
 \end{aligned}$$

where $\mathbf{H}_g^{(i)}$ is the i th component of the derivative of $\mathbf{J}_g(\mathbf{x}^*)$.

Multidimensional Newton's Method

- ▶ Newton's method corresponds to the fixed-point iteration

$$\mathbf{g}(\mathbf{x}) = \mathbf{x} - \mathbf{J}_f^{-1}(\mathbf{x})\mathbf{f}(\mathbf{x})$$

- ▶ *Note that generally Newton's method iteration has a fixed-point $\bar{\mathbf{x}}$ whenever $\mathbf{f}(\bar{\mathbf{x}}) = \mathbf{0}$, i.e. we have found a root of \mathbf{f} , namely $\mathbf{x}^* = \bar{\mathbf{x}}$.*
 - ▶ *Generally, we require that $\mathbf{J}_f(\mathbf{x}^*)$ is nonsingular, otherwise we can find nonzero solutions \mathbf{y} to $\mathbf{J}_f(\mathbf{x}^*)\mathbf{y} = \mathbf{f}(\mathbf{x}^*) = \mathbf{0}$.*
- ▶ Quadratic convergence is achieved when the Jacobian of a fixed-point iteration is zero at the solution, which is true for Newton's method:

$$\begin{aligned}\mathbf{g}(\mathbf{x}^*) &= \mathbf{x}^* - \mathbf{J}_f^{-1}(\mathbf{x}^*)\mathbf{f}(\mathbf{x}^*) \\ \mathbf{J}_g(\mathbf{x}^*) &= \mathbf{I} - \mathbf{J}_f^{-1}(\mathbf{x}^*)\mathbf{J}_f(\mathbf{x}^*) - \sum_i \underbrace{f_i(\mathbf{x}^*)}_{0} \mathbf{H}_f^{(i)}(\mathbf{x}^*) \\ &= \mathbf{I} - \mathbf{I} - \mathbf{O} = \mathbf{O}\end{aligned}$$

where $\mathbf{H}_f^{(i)}$ is the i th component of the derivative of $\mathbf{J}_f^{-1}(\mathbf{x}^*)$.

Estimating the Jacobian using Finite Differences

- ▶ To obtain $\mathbf{J}_f(\mathbf{x}_k)$ at iteration k , can use finite differences:
 - ▶ If $\mathbf{J}_f(x) \in \mathbb{R}^{m \times 1}$ (single-variate but vector-valued f), we can estimate

$$\mathbf{J}_f(\mathbf{x}_k) \approx (1/h)(\mathbf{f}(\mathbf{x}_k + h) - \mathbf{f}(\mathbf{x}_k)).$$

- ▶ More generally, the i th column of \mathbf{j}_i of the Jacobian $\mathbf{J}_f(\mathbf{x}_k)$ can be estimated by

$$\mathbf{j}_i \approx (1/h)(\mathbf{f}(\mathbf{x}_k + h\mathbf{e}_i) - \mathbf{f}(\mathbf{x}_k)).$$

- ▶ $n + 1$ function evaluations are needed: $\mathbf{f}(\mathbf{x})$ and $\mathbf{f}(\mathbf{x} + h\mathbf{e}_i)$, $\forall i \in \{1, \dots, n\}$, which correspond to $m(n + 1)$ scalar function evaluations if $\mathbf{J}_f(\mathbf{x}_k) \in \mathbb{R}^{m \times n}$.

Cost of Multivariate Newton Iteration

- ▶ What is the cost of solving $\mathbf{J}_f(\mathbf{x}_k)\mathbf{s}_k = \mathbf{f}(\mathbf{x}_k)$?
 $O(n^3)$
- ▶ What is the cost of Newton's iteration overall?
For k steps, $O(n^3k + kn^2\gamma_{\text{func-eval}})$.

Quasi-Newton Methods

In solving a nonlinear equation, seek approximate Jacobian $\mathbf{J}_f(\mathbf{x}_k)$ for each \mathbf{x}_k

- Find $\mathbf{B}_{k+1} = \mathbf{B}_k + \delta\mathbf{B}_k \approx \mathbf{J}_f(\mathbf{x}_{k+1})$, so as to approximate *secant equation*

$$\mathbf{B}_{k+1} \underbrace{(\mathbf{x}_{k+1} - \mathbf{x}_k)}_{\delta\mathbf{x}} = \underbrace{\mathbf{f}(\mathbf{x}_{k+1}) - \mathbf{f}(\mathbf{x}_k)}_{\delta\mathbf{f}}$$

Generally, the secant equation is underdetermined, as we usually need n finite-difference formulas to determine $\mathbf{J}_f(\mathbf{x}_k)$, so the secant updating methods find only approximate \mathbf{B}_{k+1} , usually as a modification of \mathbf{B}_k .

- *Broyden's method* solves the secant equation and minimizes $\|\delta\mathbf{B}_k\|_F$:

$$\delta\mathbf{B}_k = \frac{\delta\mathbf{f} - \mathbf{B}_k\delta\mathbf{x}}{\|\delta\mathbf{x}\|^2} \delta\mathbf{x}^T$$

Note that $\delta\mathbf{B}_k$ is rank-1. Consequently, we can use the Sherman-Morrison formula to update \mathbf{B}_{k+1}^{-1} with $O(n^2)$ work. Various other variants exist.

Safeguarding Methods

- ▶ Can dampen step-size to improve reliability of Newton or Broyden iteration:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{s}_k \quad \text{where} \quad \alpha_k \leq 1$$

can pick α_k so to ensure $\|\mathbf{f}(\mathbf{x}_{k+1})\| < \|\mathbf{f}(\mathbf{x}_k)\|$ or by doing a line-search to minimize $\|\mathbf{f}(\mathbf{x}_k + \alpha_k \mathbf{s}_k)\|$.

- ▶ *Trust region methods* provide general step-size control:

Establish/maintain/update region within which step is expected to be accurate. Pick each step to stay within trust region while minimizing $\|\mathbf{f}(\mathbf{x}_{k+1})\|$. Observe that the Newton-like generally seek to progress to a minima of $\|\mathbf{f}(\mathbf{x}_{k+1})\|$, and indeed much of the theory of these methods targets optimization.