

Basic NA, root finding, optimization

Motivating example: Maximum Likelihood

$x_1, \dots, x_n \sim \text{i.i.d. } N(\mu, 1)$ observations

$$L(\mu) = \prod_{i=1}^n f(x_i; \mu)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{(x_i - \mu)^2}{2}}$$

$$= \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2}$$

$$l(\mu) = -\frac{n}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2$$

The MLE of μ , $\hat{\mu}$, satisfies $\overbrace{l'(\hat{\mu}) = 0}^{\text{root finding on } l'}$, $\underbrace{\text{maximization on } l}$.

In this case $\hat{\mu}$ can be found analytically:

$$l'(\mu) = -\sum (x_i - \mu)$$

$$\Rightarrow l'(\hat{\mu}) = 0 \quad \Rightarrow \quad \hat{\mu} = \frac{1}{n} \sum x_i$$

In general, l, l' , etc cannot be evaluated analytically, particularly when multiple parameters.

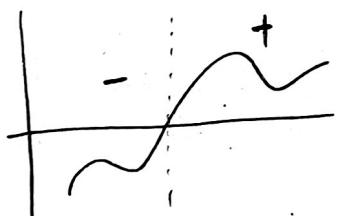
Root finding i.e. solving $f(x) = 0$
in one-dimension

Two basic ideas:

- search using values of f
- use derivative in a_0 (or approximate derivative in b_0)

Bisection

Idea: Look for sign changes of f (only $\text{sgn}(f(x))$ is used)

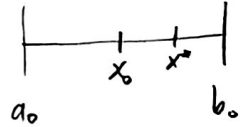


- Start with interval $[a_0, b_0]$ which contains a sign change.
- Set $x_0 = \frac{a_0 + b_0}{2}$
- If $f(a_0) f(x_0) < 0$ set $a_1 = a_0$
 $b_1 = x_0$
- If $f(x_0) f(b_0) < 0$ set $a_1 = x_0$
 $b_1 = b_0$
- Repeat to x_n compute $[a_n, b_n]$
- Take $x_n = \frac{a_n + b_n}{2}$ as estimate for root of $f(x)$.

Rate of Convergence

- If $[a_0, b_0]$ contains a sign change, then
if x^* is true root, $x_0 = \frac{a_0 + b_0}{2}$ satisfies

$$|x_0 - x^*| \leq \frac{b_0 - a_0}{2}$$



Subsequently $\underbrace{|x_n - x^*|}_{\epsilon_n} \leq \underbrace{\frac{b_0 - a_0}{2^{n+1}}}$

linear convergence

$$\frac{\epsilon_{n+1}}{\epsilon_n} = \frac{b_0 - a_0}{2^{n+2}} \frac{2^{n+1}}{b_0 - a_0} = \frac{1}{2}$$

If we require $|x_n - x^*| \leq \epsilon$, then we need

$$\frac{b_0 - a_0}{2^{n+1}} \leq \epsilon \quad = (\text{solve for } n)$$

$$\Rightarrow \log_2 \frac{b_0 - a_0}{\epsilon} \leq n+1$$

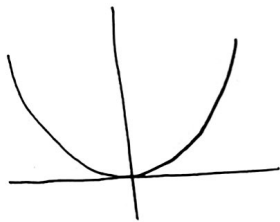
$$\Rightarrow n \geq \log_2 \frac{b_0 - a_0}{\epsilon} - 1$$

I.e. $\epsilon \approx 10^{-3} \Rightarrow n \approx 10$ if $b_0 - a_0 \approx 1$.

Main Failure Modes of Bisection

There is no failure mode! So long as $[a, b]$ contains a sign change. Bullet-proof

Note Bisection is not applicable to $f(x) = x^2$, for example.



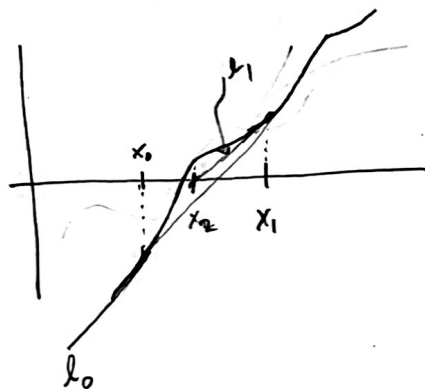
Root, but no sign change.

Newton's Method

Bisection only used function value information, how do we incorporate derivative information?

Idea Approximate f as a linear function (value + slope) and find root of the approximation - use this as estimate of true root, repeat.

Graphically:



What is the algorithm? Given initial guess x_0 ,
evaluate $f(x_0)$ and $f'(x_0)$. Form approximation near x_0 :

$$f(x) \approx f(x_0) + f'(x_0)(x-x_0) = l_0(x)$$

At the root of l_0 , call it x_1 , $l_0(x_1) = 0$.

$$\Rightarrow f(x_0) + f'(x_0)(x_1-x_0) = 0$$

$$\Rightarrow x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

and repeat to obtain:

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \quad \left. \vphantom{x_{n+1}} \right\} \text{Newton's Method}$$

Error Analysis

Assuming Newton's Method converges, (why wouldn't it converge?)

how fast does it converge?

Taylor's Thm says that (if f, f', f'' exist and are continuous near x_0)

$$f(x) = f(x_0) + f'(x_0)(x-x_0) + \frac{f''(\xi)}{2}(x-x_0)^2 \quad \text{for some } \xi \in [x_0, x]$$

\uparrow
equals

Evaluating this expression at x^* , the true root, gives :

$$(*) \quad f(x^*) = 0 = f(x_0) + f'(x_0)(x^* - x_0) + \frac{f''(\xi)}{2!} (x^* - x_0)^2$$

And Newton's Method says that

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

We are interested in the absolute error $|x^* - x_1|$

$$(**) \Rightarrow x^* - x_1 = x^* - x_0 + \frac{f(x_0)}{f'(x_0)}$$

$$\text{From Taylor, } (*) \quad x^* - x_0 = -\frac{f(x_0)}{f'(x_0)} - \frac{f''(\xi)}{2 f'(x_0)} (x^* - x_0)^2$$

Inserting into (**) we have

$$x^* - x_1 = -\frac{f''(\xi)}{2 f'(x_0)} (x^* - x_0)^2$$

Basically requires this to be bounded.

$$\Rightarrow |x^* - x_1| = \epsilon_1 = \frac{1}{2} \left| \frac{f''(\xi)}{f'(x_0)} \right| |x^* - x_0|^2 = \frac{1}{2} \left| \frac{f''(\xi)}{f'(x_0)} \right| \epsilon_0^2$$

$$\text{So } \frac{\epsilon_{n+1}}{\epsilon_n} = \frac{1}{2} \left| \frac{f''(\xi)}{f'(x_0)} \right|$$

Quadratic Convergence

Very fast, very rare - hence why Newton's Method is so powerful.

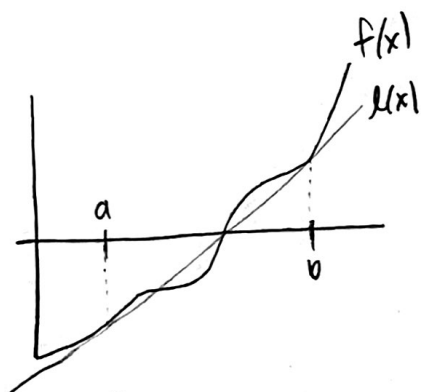
Failure Modes of Newton's Method



(1) Repeated roots

Newton's Method will "sort of" converge, but not quadratically

Example



(2) Oscillating behavior

Happens rarely, but happens because x_0 is not close enough to root for Newton to converge.

Quick programming demo

Quasi-Newton Methods

Newton: $x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$

→ Replace with an approximation, e.g.

$$f'(x_n) \approx \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}} \quad \} \Rightarrow \text{Secant Method}$$

Multivariate Newton's Method

Newton's Method can be generalized to systems of non-linear equations:

$$f_1(x_1, x_2, \dots, x_n) = 0$$

$$\vdots$$

$$f_n(x_1, \dots, x_n) = 0$$

$$\Leftrightarrow \vec{f}(\vec{x}) = \vec{0}.$$

Vector notation.

An analogous Taylor approximation can be made:

$$(***) \quad \vec{f}(\vec{x}) \approx \vec{f}(\vec{x}_0) + \vec{J}(\vec{x}_0) (\vec{x} - \vec{x}_0)$$

where \vec{J} is the Jacobian matrix:

$$J_{ij}(\vec{x}_0) = \frac{\partial f_i}{\partial x_j}(\vec{x}_0)$$

Assuming $\vec{J}^{-1}(\vec{x}_0)$ exists, the multivariate Newton's method

is given by: (by evaluating (***) at true root \vec{x}^*)

$$\vec{x}_{n+1} = \vec{x}_n - \underbrace{\vec{J}^{-1}(\vec{x}_0)}_{\substack{n \times n \\ \text{matrix}}} \underbrace{\vec{f}(\vec{x}_0)}_{\substack{n \times 1 \\ \text{vector}}}$$

Convergence can also be shown to be quadratic:

$$\|\vec{x}_{n+1} - \vec{x}^*\| \approx C \|\vec{x}_n - \vec{x}^*\|^2,$$

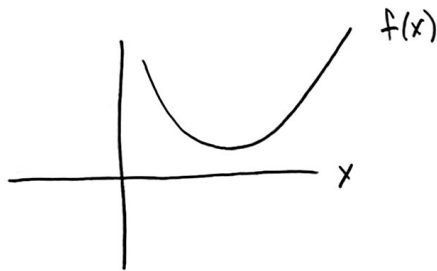
Optimization

There are two basic ways to classify optimization problems:

convex vs. non-convex

unconstrained vs. constrained

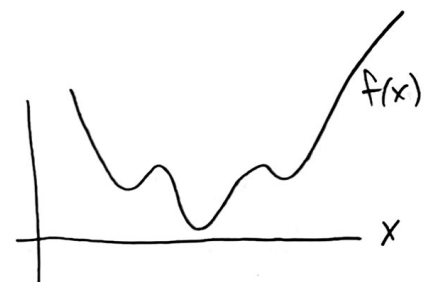
In one-dimension:



convex

EASY

$f'(x) = 0$ has one unique solution



non-convex

HARD hard to design algorithms that do not get "stuck" in local minima

Eg. $f'(x) = 0$ has many solutions

Unconstrained

$$\max_{x \in \mathbb{R}} f(x)$$

Constrained

$$\begin{aligned} &\max_x f(x) \\ &\text{subject to } g(x) = 0 \\ &\text{or } \leq \end{aligned}$$

} Example
 $\max_x e^x$
such that $|x-5| \leq 0$

We will only talk about methods for

convex unconstrained problems.

(See dedicated optimization course at Courant for max.)

Comparison with Root finding

$\max_{\vec{x}} f(\vec{x})$ ← notice f is scalar - valued

At the maximum, $\nabla f(\vec{x}^*) = \vec{0}$ — so we must find roots of $\nabla f = \vec{0}$.

Simplest method: Gradient Descent (Ascent)

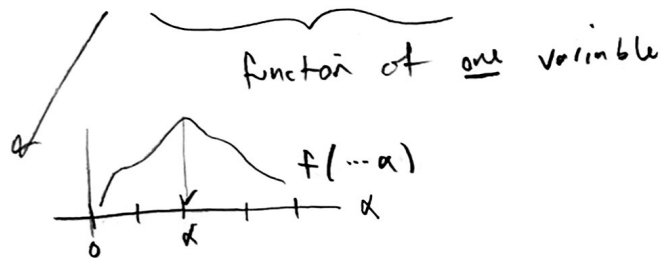
∇f points in the direction of steepest change of the function f — starting at \vec{x}_0 , move in this direction to obtain next estimate of maximum:

$$\vec{x}_1 = \vec{x}_0 + \alpha \nabla f(\vec{x}_0)$$

↑ Choose α so that $f(\vec{x}_1) \geq f(\vec{x}_0)$

Most elementary method: line search

choose α to ^{approximately} maximize $f(\vec{x}_0 + \alpha \nabla f(\vec{x}_0))$



For convex problems, and proper choice of α ,

this method is linearly convergent (like bisection).

Values of ∇f are used to solve $\nabla f = \vec{0}$, but not derivatives of ∇f .

(Foreshadow to stochastic gradient descent)

A Newton Method

Newton's method generalizes straightforwardly to optimization, except we need 2nd partial derivatives of f .

Recall, we are trying to find a zero of ∇f , so approx by a Taylor series:

$$\nabla f(\vec{x}) \approx \nabla f(\vec{x}_0) + H(\vec{x}_0)(\vec{x} - \vec{x}_0)$$

where $\nabla f = \begin{pmatrix} \partial f / \partial x_1 \\ \partial f / \partial x_2 \\ \vdots \\ \partial f / \partial x_n \end{pmatrix}$

and $H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$

this is known as

the Hessian matrix: symmetric
↓
what else?

⇒ the Newton iteration is then

$$\vec{x}_{n+1} = \vec{x}_n - H^{-1}(\vec{x}_n) \nabla f(\vec{x}_n) \quad] \text{ computational complexity?}$$

Often H cannot be computed, (for whatever reason), or it is too expensive to compute. Approximating it by \tilde{H} gives us... a "Quasi-Newton method".

→ Mention step-size control.

→ Mention General form of optimization algorithms: $\vec{x}_{k+1} = \vec{x}_k + \alpha_k \vec{p}_k$