

Topic

Continuous Optimization II

Today

Applications of gradient descent

- Linear system solving

Newton's method (another iterative method for continuous optimization)

Applications:

- Computing roots
- Unconstrained minimization

Review of last time

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(x) \\ \vdots \\ \frac{\partial f}{\partial x_n}(x) \end{bmatrix} \in \mathbb{R}^n$$

Gradient Descent Method

start with x^0

Repeat

$$x^{k+1} = x^k - \eta_k \nabla f(x^k)$$

step size can vary with k

based on Taylor Expansion

$$f(x) = \underbrace{f(x^i) + \nabla f(x^i)(x - x^i)}_{\text{linear approximation of } f} + \underbrace{\frac{1}{2}(x - x^i)^T \nabla^2 f(x^i)(x - x^i) + \dots}_{\text{"error" hopefully small}}$$

$\approx \eta_k$
 \downarrow

Known fact

gradient descent converges to \tilde{x}

s.t. $\nabla f(\tilde{x}) = 0$

\Uparrow
 might not even be local min:
 could be local max or saddle pt.
 but all local mins are critical!
 (if convex, all critical pts are local mins)

Linear Systems Solving

Given system of linear equations

$$a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1$$

$$a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2$$

⋮

$$a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n = b_m$$

Find x_1, \dots, x_n satisfying equations

OR equivalently :

Given $A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \in \mathbb{R}^{m \times n}$ $b = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix} \in \mathbb{R}^m$

Find $x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n$

such that $Ax = b$

Today $m = n$ (A square matrix)
 A invertible (soln exists)

How to solve?

1st idea:

invert A :

compute A^{-1} &

$$\text{set } x^* = A^{-1}b$$

← natural approach!

why not?

- can be slow (fastest algorithm uses $\theta(n^{2.373})$ time)
- computing A^{-1} can involve very large numbers since might need to divide by small numbers
- if A is sparse (few non zero entries) then A^{-1} might not be

but can we hope to improve?

2nd idea:

- phrase as unconstrained minimization problem
- use gradient descent!

Consider

$$f(x) = \frac{1}{2} x^T A x - b x \quad \Leftarrow \text{Quadratic function}$$

Why did we pick this crazy choice of f ?
(What does it have to do with our goal?)

What is ∇f ?

this
is looking
more
promising!!

\Rightarrow

$$\nabla f(x) = Ax - b$$

$$\text{Note: } \nabla b x = b$$

$$\nabla (x^T B x) = 2 B x$$

What does gradient descent do?

it finds an extremum - where $\nabla f(\tilde{x}) \approx 0$

if we find such an \tilde{x} we have

$$0 \approx \nabla f(\tilde{x}) = A \tilde{x} - b$$

$$\text{so } A \tilde{x} \approx b$$

Comment: in fact, we picked f to be integral of equation we wanted to solve a solution!!

Advantages:

- Fast
- no division

Root finding

given $f: \mathbb{R}^n \rightarrow \mathbb{R}$

find x^* s.t. $f(x^*) = 0$

\Leftarrow very different goal, right?

actually, we'll see - that it is very related

e.g. $f(x) = x^2 - 2 \Rightarrow x^* = \pm\sqrt{2}$

idea 1 use gradient descent again

define $g(x) = \frac{x^3}{3} - 2x$

\Leftarrow integral of f

$$g'(x) = \frac{d}{dx} \left(\frac{x^3}{3} - 2x \right) = x^2 - 2$$

setting $g'(x) = x^2 - 2 = 0$

gradient descent

\Rightarrow solution $\tilde{x} = \pm\sqrt{2}$

idea 2 Newton's method (also Newton-Raphson)

(also iterative approach!)

As before approx $f(x)$ using Taylor series:

$$f(x) = f(x^k) + f'(x^k)(x - x^k) + \underbrace{\frac{1}{2} f''(x^k)(x - x^k)^2 + \dots}_{\text{error}}$$

\hat{f}_k is linear approx of f around x^k

since want x s.t. $f(x) = 0$:

$$0 = \hat{f}_k(x^{k+1}) = f(x^k) + f'(x^k)(x^{k+1} - x^k)$$

$$\Rightarrow x^{k+1} = x^k - \frac{f(x^k)}{f'(x^k)}$$

\Leftarrow set $\hat{f}_k(x^{k+1})$ to 0 & solve for x^{k+1}

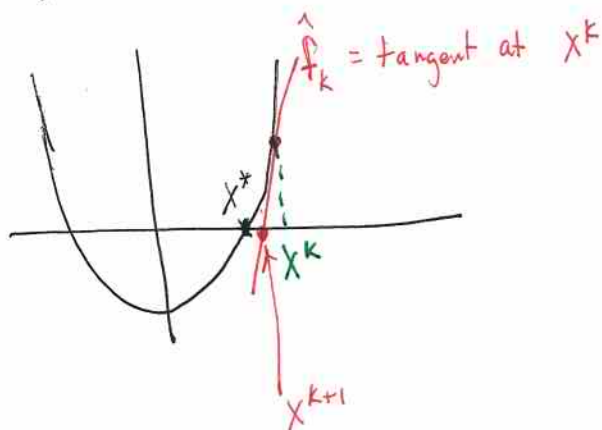
(6)

Newton's method for $n=1$ (Newton-Raphson)

start with x^0

repeat

$$x^{k+1} = x^k - \frac{f(x^k)}{f'(x^k)}$$



General case $n \geq 1$

$$f(x) = \underbrace{f(x^k) + \nabla f(x^k)^T (x - x^k)}_{\hat{f}_k(x)} + \underbrace{\frac{1}{2} (x - x^k)^T \nabla^2 f(x^k) (x - x^k) + \dots}_{\text{error}}$$

Setting

$$0 = \hat{f}_k(x) = f(x^k) + \nabla f(x^k)^T (x - x^k)$$

$$\text{get } x^{k+1} = x^k - \frac{f(x^k)}{\|\nabla f(x^k)\|^2} \nabla f(x^k)$$

$= \frac{f(x^k)}{f'(x^k)} \text{ when } n=1$

note

f is non-linear & we are ignoring that
step size might not be small enough to make error small

\Rightarrow convergence not as robust as for GD

i.e. might not converge

but if it does converge, then it is

Turns out that error is being squared! **FAST**
So if error > 1 , big problem
but if error < 1 , converges fast

(7)

Newton's method for ^{square} roots:

$$f(x) = x^2 - a$$

$$x^{(k+1)} = x^{(k)} - \left(\frac{\overset{f(x^{(k)})}{x^{(k)2} - a}}{\underset{f'(x^{(k)})}{2x^{(k)}}} \right) = x^{(k)} - \frac{x^{(k)}}{2} + \frac{a}{2x^{(k)}} = \frac{1}{2} \left[\underbrace{x^{(k)} + \frac{a}{x^{(k)}}}_{\substack{\text{average} \\ \text{of } x^{(k)} \\ + \frac{a}{x^{(k)}}}} \right]$$

↑↑
dates back
to Babylonians

Example $a=2$

$$x^{(0)} = \underline{1.00000\dots}$$

$$x^{(1)} = \underline{1.50000\dots}$$

$$(1 + \frac{2}{1})/2 = 3/2$$

$$x^{(2)} = \underline{\underline{1.41666\dots}}$$

$$(\frac{3}{2} + \frac{2}{3/2})/2 = 17/12$$

$$x^{(3)} = \underline{\underline{1.414215685}}$$

$$577/408$$

$$x^{(4)} = \underline{\underline{1.414213562}}$$

correct digits nearly doubles !!
(quadratic convergence)

Error analysis

⑧

Assume multiplicative error at stage k is $(1 + \epsilon_k)$

$$\text{i.e. } x^{(k)} = \sqrt{a} (1 + \epsilon_k)$$

$$\text{then } x^{(k+1)} = \frac{x^{(k)} + \left(\frac{a}{x^{(k)}}\right)}{2}$$

$$= \frac{\sqrt{a} (1 + \epsilon_k) + \left(\frac{a}{\sqrt{a} (1 + \epsilon_k)}\right)}{2}$$

$$= \sqrt{a} \left[\frac{1 + \epsilon_k + \frac{1}{1 + \epsilon_k}}{2} \right] = \sqrt{a} \left[\frac{1 + 2\epsilon_k + \epsilon_k^2 + 1}{2(1 + \epsilon_k)} \right]$$

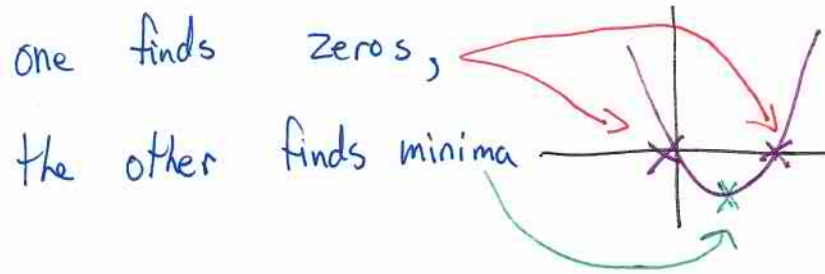
$$= \sqrt{a} \left[\frac{2 + 2\epsilon_k + \epsilon_k^2}{2(1 + \epsilon_k)} \right]$$

$$= \sqrt{a} \left[1 + \underbrace{\frac{\epsilon_k^2}{2(1 + \epsilon_k)}}_{= \epsilon_{k+1}} \right]$$

\Leftarrow quadratic
convergence
for $\epsilon_0 < 1$

(9)

Newton's method & unconstrained minimization are different!



But Newton's method & unconstrained minimization are also very related!

- we already saw that minimization of "integral" can be used to find zeroes
- finding zeroes can be used to solve minimization too!

e.g. to find $x^* = \arg\min_x f(x)$ for convex f

compute root of $g(x) = \|\nabla f(x)\|^2$

$$g(x^*) = 0 \Rightarrow \nabla f(x^*) = 0 \Rightarrow x^* = \arg\min_x f(x)$$

\uparrow
 convexity