

CSC338 Lecture 11

Computational Problem: Unconstrained optimization
 $f: \mathbb{R} \rightarrow \mathbb{R}$, find a local minimum of $f(x)$.

① Golden Section Search: If f is unimodal on $[a, b]$ then we can iteratively shrink the interval in which x^* lies in. Only requires evaluation of f , not derivatives. However, requires the unimodal property. Always converges.

② Newton's Method: Idea - approximate $f(x)$ using a Quadratic function, and find a critical point of the approximation.

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

Requires derivative and second derivative.

Does not always converge. Does not always converge to a min.

Today $f: \mathbb{R}^n \rightarrow \mathbb{R}$, find local minimum of $f(x)$.

- 1) Newton's method
- 2) Gradient Descent

Newton's method for $f: \mathbb{R}^n \rightarrow \mathbb{R}$

Recall, for $f: \mathbb{R} \rightarrow \mathbb{R}$, we have the Taylor Series Expansion

$$f(x+h) = f(x) + f'(x)h + \frac{1}{2}f''(x)h^2 + O(h^3)$$

This result extends to $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$f(\underline{x} + \underline{s}) = f(\underline{x}) + \underbrace{\nabla f(\underline{x})^T}_{\text{Gradient}} \underbrace{\underline{s}}_{n \times 1} + \frac{1}{2} \underbrace{\underline{s}^T}_{1 \times n} \underbrace{H_f(\underline{x})}_{n \times n} \underbrace{\underline{s}}_{n \times 1}$$

$$\begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

Idea: In each iteration, we have an estimate \underline{x}_k of a minimum approximate $f(\underline{x})$ with

$$f_k(\underline{x}) = f(\underline{x}_k) + \nabla f(\underline{x}_k)^T \underline{s} + \frac{1}{2} \underline{s}^T H_f(\underline{x}_k) \underline{s}$$

\parallel
 $(\underline{x}_k + \underline{s})$

And find \underline{s} that (locally) minimizes $f_k(\underline{x}_k + \underline{s})$.

Critical points of f_k has:

$$\nabla_{\underline{s}} f_k(\underline{x}_k + \underline{s}) = 0 \quad (\text{like } f'(x) = 0)$$

$$\Rightarrow \underbrace{H_f(\underline{x}_k)}_{n \times n} \underbrace{\underline{s}_k}_{n \times 1} = - \underbrace{\nabla f(\underline{x}_k)}_{n \times 1} \quad (\text{like } f''(x)h = -f'(x))$$

oblem

$f''(x^*)$
large

δ small

ϵ is small

so

why

$f'(x) = 0$

$\mathbb{R} \rightarrow \mathbb{R}$

x_k

$\mathbb{R}^{n \times n}$

2x1
4

Newton's method update rule:

$$\underline{x}_{k+1} = \underline{x}_k + \underline{s}_k \quad \left(\underline{s}_k = -(\underline{H}_f(\underline{x}_k))^{-1} \nabla f(\underline{x}_k) \right)$$

One major disadvantage of Newton's method is that computing the Hessian $\underline{H}_f(\underline{x})$ is expensive. $\underline{H}_f(\underline{x}) \in \mathbb{R}^{n \times n}$

eg // $f(\underline{x}) = x_1^4 + x_1^2 x_2 + x_1^2 + 2x_2^2 + x_2$

We wish to find a min, starting with $\underline{x}_0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$$\nabla f(\underline{x}) = \begin{bmatrix} 4x_1^3 + 2x_1x_2 + 2x_1 \\ x_1^2 + 4x_2 + 1 \end{bmatrix}$$

$$\underline{H}_f(\underline{x}) = \begin{bmatrix} 12x_1 + 2x_2 + 2 & 2x_1 \\ 2x_1 & 4 \end{bmatrix}$$

$$\nabla f(\underline{x}_0) = \begin{bmatrix} 8 \\ 6 \end{bmatrix}$$

$$\underline{H}_k(\underline{x}_0) = \begin{bmatrix} 16 & 2 \\ 2 & 4 \end{bmatrix}$$

Solve $\underline{s}_0 = \begin{bmatrix} -2/5 \\ -4/5 \end{bmatrix}$

Reading contour plots:

Each ~~of~~ ^A the curves in the contour plot shows points with the ~~same~~ same $f(x)$ value.

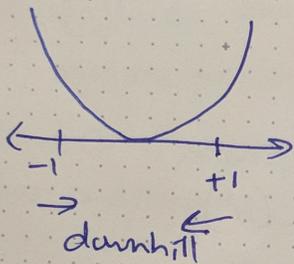
Steepest descent (gradient descent)

We'll discuss a variation of the algorithm from your textbook. This version is typically used in machine learning.

Key idea: ~~the negative gradient~~

The gradient of a differentiable function points uphill
: negative gradient : : : downhill
Toward x with lower values of $f(x)$.

Example in 1D: $f(x) = x^2 : \mathbb{R} \rightarrow \mathbb{R}$



$$f'(-1) = 2(-1) = -2$$

\Rightarrow negative gradient points to the right

$$f'(1) = 2(1) = 2$$

\Rightarrow negative gradient points to the left

Example in 2D: $f(x) = x_1^4 + x_1 x_2 + x_1^2 + 2x_2^2 + x_2$

At $x = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ $\nabla f(x) = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ so downhill is in the direction $\begin{bmatrix} 0 \\ -1 \end{bmatrix}$

$x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\nabla f(x) = \begin{bmatrix} 8 \\ 6 \end{bmatrix}$ so : : : : : $\begin{bmatrix} -8 \\ -6 \end{bmatrix}$

Note: The gradient is always perpendicular to the contour!

Why does this work?

Using Taylor's Thm, we can locally approximate f using

$$f(x+s) = f(x) + \nabla f(x)^T s$$

If $\nabla f(x) \neq 0$, then if we let $s = -\nabla f(x) \cdot c$ for some c small

Then $f(x+s) = f(x) +$ (something negative) move x towards $-\nabla f(x)$

It turns out that $-\nabla f(x)$ is, locally, the direction of the steepest descent.

Steepest Descent for finding a minima of $f: \mathbb{R}^n \rightarrow \mathbb{R}$

Start from initial guess \underline{x}_0 , and update

$$\underline{x}_{k+1} = \underline{x}_k - \alpha_k \nabla f(\underline{x}_k) \quad \alpha_k \in \mathbb{R}, \alpha_k \geq 0$$

How to find α_k ?
(or choose)

Textbook: use of line search: find α_k ~~so that~~ ^{to minimize}

$$g(\alpha_k) = f(\underline{x}_k - \alpha_k \nabla f(\underline{x}_k)) : \mathbb{R} \rightarrow \mathbb{R}$$

Gradient

Descent: Set α_k to a constant α

called the learning rate

Steepest descent is usually reliable: can always make progress provided that the gradient is nonzero.

However, this method is myopic in its view of function's ~~behavior~~ ^{behavior}

only take into account local information

eg // $f(x) = 0.5x_1^2 + 2.5x_2^2$

$$\nabla f(x) = \begin{bmatrix} x_1 \\ 5x_2 \end{bmatrix}$$

... slow slide 29 & 30.

Comparison:

Steepest Descent: $\underline{x}_{k+1} = \underline{x}_k - \alpha_k \nabla f(\underline{x}_k)$

Newton: $\underline{x}_{k+1} = \underline{x}_k - H_f(\underline{x}_k)^{-1} \nabla f(\underline{x}_k)$

Quasi-Newton: $\underline{x}_{k+1} = \underline{x}_k - \alpha_k B_k^{-1} \nabla f(\underline{x}_k)$

like assuming $H_f(\underline{x}_k) = \alpha_k I$

Where B_k is an approximation of the Hessian matrix. (eg like secant method)

\Rightarrow BFGS.

Revisit hw grade prediction: predict HW3 given HW1, HW2 grades.

$$A = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} \\ a_{11}^{(2)} & a_{12}^{(2)} \\ \vdots & \vdots \\ a_{11}^{(73)} & a_{12}^{(73)} \end{bmatrix}$$

column = one hw

$$\underline{b} = \begin{bmatrix} b_{11} \\ b_{12} \\ \vdots \\ b_{73} \end{bmatrix}$$

row = one student

Problem find $\underline{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ to minimize $\|A\underline{x} - \underline{b}\|_2$

We can treat this as a non-linear optimization problem:
(see slides for rest of problem) ~~(see next two pages)~~