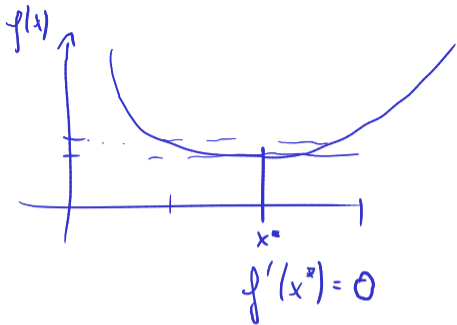


- Changes:

- no AK office hours today ✓
 - very short AK office hours Tue ✓
 - no AK office hours Thu / Luke teaching
 - no class Nov 8 ✓
- 4-credit assignment 1 posted
- Examlet 4 ongoing

$$Ax = u(v^T x)$$
$$A = \begin{array}{|c|} \hline \square \\ \hline \end{array} = u v^T \quad \begin{array}{|c|} \hline \square \\ \hline \end{array} = \begin{array}{|c|} \hline \square \\ \hline \end{array}$$

~~$= u v^T$~~

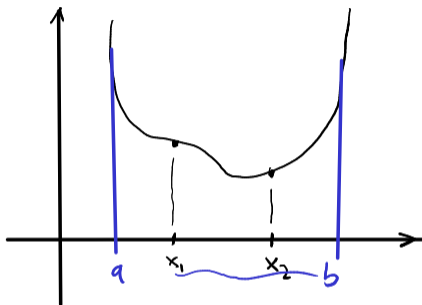


unimodal: $f(x^*+h) = f(x^*) + f''(x^*) \frac{h^2}{2}$



Golden Section Search

Suppose we have an interval with f unimodal:



Would like to maintain unimodality.

Pick x_1, x_2

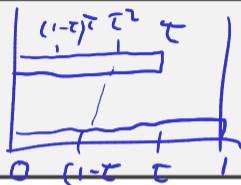
If $f(x_1) > f(x_2)$, reduce to (x_1, b)

If $f(x_1) \leq f(x_2)$, reduce to (a, x_2)

Golden Section Search: Efficiency

Where to put x_1, x_2 ?

$$x_1 = a + (1-\tau)(b-a)$$
$$x_2 = a + \tau(b-a)$$



Demo: Golden Section Proportions [cleared]

$$\tau^2 = 1 - \tau$$

$$\tau = (\sqrt{5} - 1) / 2$$

"golden section search"


Convergence rate?

linear

Newton's Method

$$\text{solve } f'(x) = 0$$

Reuse the Taylor approximation idea, but for optimization.


$$f(x+h) \approx f(x) + f'(x) \cdot h + \frac{f''(x)}{2} \cdot h^2 = \hat{f}(h)$$
$$\hat{f}'(h) = f'(x) + f''(x) \cdot h \stackrel{!}{=} 0$$
$$h = - \frac{f'(x)}{f''(x)}$$

$x_0 =$ (starting guess)

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

Demo: Newton's Method in 1D [cleared]

Steepest Descent/Gradient Descent

Given a scalar function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at a point \mathbf{x} , which way is down?

Direction of steepest descent: $-\nabla f$



- $\vec{x}_0 =$ (starting guess)
- $\vec{s}_n = -\nabla f(\vec{x}_n)$
- Use a line search to choose α_n to minimize the function
 $\alpha \mapsto f(\vec{x}_n + \alpha \vec{s}_n)$
- $\vec{x}_{n+1} = \vec{x}_n + \alpha_n \vec{s}_n$

Demo: Steepest Descent [cleared] (Part 1)

Steepest Descent: Convergence

Consider quadratic model problem:

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{c}^T \mathbf{x}$$

where A is SPD. (A good model of f near a minimum.)

$$\mathbf{x}^T \mathbf{A} \mathbf{x} \rightarrow \begin{pmatrix} 0.5 & 0 \\ 0 & 2.5 \end{pmatrix}$$

Define error $\vec{e}_k = \vec{x}_k - \vec{x}^*$

$$\|\vec{e}_{k+1}\|_A = \sqrt{\vec{e}_{k+1}^T \mathbf{A} \vec{e}_{k+1}}$$

$$\leq \frac{\sigma_{\max}(\mathbf{A}) - \sigma_{\min}(\mathbf{A})}{\sigma_{\max}(\mathbf{A}) + \sigma_{\min}(\mathbf{A})} \|\vec{e}_k\|_A$$


$$= \frac{\kappa(\mathbf{A}) - 1}{\kappa(\mathbf{A}) + 1} \cdot \|\vec{e}_k\|_A$$

$$\|\mathbf{x}\|_A = \sqrt{\mathbf{x}^T \mathbf{A} \mathbf{x}}$$

(only SPD)

Hacking Steepest Descent for Better Convergence

Extrapolation methods:


$$x_{n+1} = x_n - \alpha_n \nabla f(x_n) + \beta_n (\vec{x}_n - \vec{x}_{n-1})$$

Heavy ball method:

$$\alpha_n = \alpha \quad \beta_n = \beta$$

Demo: Steepest Descent [cleared] (Part 2)

Optimization in Machine Learning

What is *stochastic gradient descent (SGD)*?

$$f(\vec{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\vec{x})$$

"batch"

$$\vec{x}_{t+1} = \vec{x}_t - \alpha \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_{it})$$

"minibatch"

↑
same part of this sum

"ADAM"

moving averages of ∇ and the square of the gradient

Conjugate Gradient Methods

Can we optimize in *the space spanned* by the last two step directions?

$$(\alpha_a, \beta_a) = \operatorname{argmin}_{\alpha, \beta} \left[f(x_a - \alpha \nabla f(x_a) + \beta (x_a - x_{a-1})) \right]$$

Demo: Conjugate Gradient Method [cleared]

Nelder-Mead Method



Idea:



Demo: Nelder-Mead Method [cleared]