

# CS 450: Numerical Analysis<sup>1</sup>

## Numerical Optimization

University of Illinois at Urbana-Champaign

---

<sup>1</sup>*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](#)).*

# Numerical Optimization

- ▶ Our focus will be on *continuous* rather than *combinatorial* optimization:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to} \quad \mathbf{g}(\mathbf{x}) = \mathbf{0} \quad \text{and} \quad \mathbf{h}(\mathbf{x}) \leq \mathbf{0}$$

- ▶ We consider linear, quadratic, and general nonlinear optimization problems:





## Optimality Conditions

- ▶ If  $\mathbf{x}$  is an interior point in the feasible domain and is a local minima,

$$\nabla f(\mathbf{x}) = \left[ \frac{df}{dx_1}(\mathbf{x}) \quad \cdots \quad \frac{df}{dx_n}(\mathbf{x}) \right]^T = \mathbf{0} :$$

- ▶ *Critical points*  $\mathbf{x}$  satisfy  $\nabla f(\mathbf{x}) = \mathbf{0}$  and can be minima, maxima, or saddle points:

## Hessian Matrix

- ▶ To ascertain whether a critical point  $\boldsymbol{x}$ , for which  $\nabla f(\boldsymbol{x}) = \mathbf{0}$ , is a local minima, consider the *Hessian matrix*:
  
- ▶ If  $\boldsymbol{x}^*$  is a minima of  $f$ , then  $\boldsymbol{H}_f(\boldsymbol{x}^*)$  is positive semi-definite:

## Optimality on Feasible Region Border

- ▶ Given an equality constraint  $\mathbf{g}(\mathbf{x}) = \mathbf{0}$ , it is no longer necessarily the case that  $\nabla f(\mathbf{x}^*) = \mathbf{0}$ . Instead, it may be that directions in which the gradient decreases lead to points outside the feasible region:

$$\exists \boldsymbol{\lambda} \in \mathbb{R}^n, \quad -\nabla f(\mathbf{x}^*) = \mathbf{J}_g^T(\mathbf{x}^*)\boldsymbol{\lambda}$$

- ▶ Such *constrained minima* are critical points of the Lagrangian function  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})$ , so they satisfy:

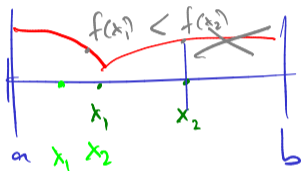
$$\nabla \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}) = \begin{bmatrix} \nabla f(\mathbf{x}^*) + \mathbf{J}_g^T(\mathbf{x}^*)\boldsymbol{\lambda} \\ \mathbf{g}(\mathbf{x}^*) \end{bmatrix} = \mathbf{0}$$





# Golden Section Search

- Given bracket  $[a, b]$  with a unique minimum ( $f$  is *unimodal* on the interval), *golden section search* considers points  $f(x_1), f(x_2)$ ,  $a < x_1 < x_2 < b$  and discards subinterval  $[a, x_1]$  or  $[x_2, b]$ :



- Since one point remains in the interval, golden section search selects  $x_1$  and  $x_2$  so one of them can be effectively reused in the next iteration:

$$\frac{(\sqrt{5}-1)}{2} (b-a) \quad \text{is distance from } a \text{ to } x_2$$

$$\frac{(\sqrt{5}-1)}{2} (b-a) \quad \text{is distance from } x_1 \text{ to } b$$








# General Multidimensional Optimization

- ▶ Direct search methods by simplex (*Nelder-Mead*):

- ▶ Steepest descent: find the minimizer in the direction of the negative gradient:

$$\alpha_k = \min_{\alpha} f(x_k - \alpha \nabla f(x_k)) \quad \left| \quad x_{k+1} = x_k - \alpha_k \nabla f(x_k) \right.$$


via Golden Section  
or any 1D optimization scheme

# Convergence of Steepest Descent

- ▶ Steepest descent converges linearly with a constant that can be arbitrarily close to 1:

in worst case  $\alpha_k \approx 0$ , always  $\alpha_k > 0$

how fast  $\nabla f$  changes guides convergence

- ▶ Given quadratic optimization problem  $f(x) = \frac{1}{2}x^T Ax + c^T x$  where  $A$  is symmetric positive definite, the error  $e_k = x_k - x^*$  satisfies

$$\|e_k\|_A = \frac{\sigma_{\max}(A) - \sigma_{\min}(A)}{\sigma_{\max}(A) + \sigma_{\min}(A)} \|e_{k-1}\|_A$$

$\frac{\kappa(A) - 1}{\kappa(A) + 1}$

$e_k^T A e_k$ , only norm if  $A$  is SPD

$$h(x_k + s) = h(x_k) + \underbrace{\nabla h(x_k)}_c s + \frac{1}{2} s^T H_h(x_k) s + \dots$$

$$\nabla f(x) = Ax + c$$

$$0 = Ax + c$$

optimality condition

min f(x)

quadratic

linear term

## Gradient Methods with Extrapolation

- ▶ We can improve the constant in the linear rate of convergence of steepest descent by leveraging *extrapolation methods*, which consider two previous iterates (maintain *momentum* in the direction  $x_k - x_{k-1}$ ):

$$x_{k+1} = f \left( x_k - \underbrace{\alpha_k}_{\text{parameter}} \underbrace{\nabla f(x_k)}_{\text{gradient}} + \underbrace{\beta_k (x_k - x_{k-1})}_{\text{momentum}} \right)$$

- ▶ The *heavy ball method*, which uses constant  $\alpha_k = \alpha$  and  $\beta_k = \beta$ , achieves better convergence than steepest descent:

$$\|e_k\|_A = \frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \|e_{k-1}\|_A \quad \left. \vphantom{\|e_k\|_A} \right\} \begin{array}{l} \text{near-optimal} \\ \text{same as Conjugate} \\ \text{Gradient} \end{array}$$

# Conjugate Gradient Method

- ▶ The *conjugate gradient method* is capable of making the optimal choice of  $\alpha_k$  and  $\beta_k$  at each iteration of an extrapolation method:

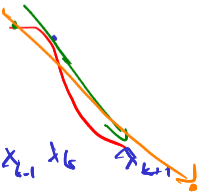
$$x_{k+1} = \underset{x}{\operatorname{argmin}} \left( f \left( x_k - \underbrace{\alpha_k \nabla f(x_k)}_{\text{gradient}} + \underbrace{\beta_k (x_k - x_{k-1})}_{\text{momentum}} \right) \right)$$

optimal  $\alpha_k, \beta_k$  at each step, Lanczos implicitly

- ▶ Parallel tangents implementation of the method proceeds as follows

- steepest descent from  $x_k$  to generate  $\hat{x}_k$
- minimize along the line from  $x_{k-1}$  to  $\hat{x}_k$  to produce  $x_{k+1}$

CG converges in  $n$  steps





## Conjugate Gradient as a Krylov Subspace Method

- ▶ Conjugate Gradient finds the minimizer of  $f(x) = \frac{1}{2}x^T Ax + c^T x$  within the Krylov subspace of  $A$ :

$$x^* = \operatorname{argmin}_{x \in \mathcal{K}_k(A, c)}$$

$$f(x) = \operatorname{argmin}_{\substack{x = Q_k y \\ y \in \mathbb{R}^k}} \underbrace{\frac{1}{2} y^T Q_k^T A Q_k y + c^T Q_k y}_{\frac{1}{2} y^T T_k y + \underbrace{c^T Q_k}_{\|c\|_2 e_1} y}$$

$$\nabla f(y) = T_k y + \|c\|_2 e_1 = 0$$

$$y = -\|c\|_2 T_k^{-1} e_1$$

$$x_k^* = Q_k y = -\|c\|_2 Q_k T_k^{-1} e_1$$

## Newton's Method

- ▶ Newton's method in  $n$  dimensions is given by finding minima of  $n$ -dimensional quadratic approximation:

$$f(x_k + s) \approx \hat{f}(s) = f(x_k) + s^T \nabla f(x_k) + \frac{1}{2} s^T H_f(x_k) s \quad \text{Taylor expansion}$$

$$\nabla \hat{f}(s) = 0 : 0 = \nabla f(x_k) + H_f(x_k) s$$

$$s = -H_f^{-1}(x_k) \nabla f(x_k)$$

$$x_{k+1} = x_k + s$$

$$\nabla f(x) = 0$$

equivalent Newton  
quadratic convergence



## Nonlinear Least Squares

- ▶ An important special case of multidimensional optimization is *nonlinear least squares*, the problem of fitting a nonlinear function  $f_{\mathbf{x}}(t)$  so that  $f_{\mathbf{x}}(t_i) \approx y_i$ :
  
- ▶ We can cast nonlinear least squares as an optimization problem and solve it by Newton's method:





## Constrained Optimization Problems

- ▶ We now return to the general case of *constrained* optimization problems:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to} \quad \mathbf{g}(\mathbf{x}) = \mathbf{0} \quad \text{and} \quad \mathbf{h}(\mathbf{x}) \leq \mathbf{0}$$

- ▶ Generally, we will seek to reduce constrained optimization problems to a series of unconstrained optimization problems:

## Lagrangian Duality

- ▶ The Lagrangian function with constraints  $\mathbf{g}(\mathbf{x}) = \mathbf{0}$  and  $\mathbf{h}(\mathbf{x}) \leq \mathbf{0}$  is

- ▶ The Lagrangian dual problem is an unconstrained optimization problem:

$$\max_{\boldsymbol{\lambda}} q(\boldsymbol{\lambda}), \quad q(\boldsymbol{\lambda}) = \begin{cases} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) & \text{if } \boldsymbol{\lambda} \geq \mathbf{0} \\ -\infty & \text{otherwise} \end{cases}$$

The unconstrained optimality condition  $\nabla q(\boldsymbol{\lambda}^*) = \mathbf{0}$ , implies



## Sequential Quadratic Programming

- ▶ *Sequential quadratic programming (SQP)* reduces a nonlinear equality constrained problem to a sequence of constrained quadratic programs via a Taylor expansion of the Lagrangian function  $\mathcal{L}_f(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})$ :
  
- ▶ SQP ignores the constant term  $\mathcal{L}_f(\mathbf{x}_k, \boldsymbol{\lambda}_k)$  and minimizes  $s$  while treating  $\delta$  as a Lagrange multiplier:

## Sequential Quadratic Programming

- ▶ From a different viewpoint, sequential quadratic programming corresponds to using Newton's method to solve the nonlinear equations,

$$\nabla \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \begin{bmatrix} \nabla f(\mathbf{x}) + \mathbf{J}_g^T(\mathbf{x})\boldsymbol{\lambda} \\ \mathbf{g}(\mathbf{x}) \end{bmatrix} = \mathbf{0}$$

## Active Set Methods

- ▶ To use SQP for an inequality constrained optimization problem, consider at each iteration an *active set* of constraints:
  
- ▶ The Karush-Kuhn-Tucker (KKT) optimality conditions given the generalized Lagrangian function  $\mathcal{L}(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\nu}) = f(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}) + \boldsymbol{\nu}^T \mathbf{h}(\mathbf{x})$  are

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{0}$$

$$\mathbf{g}(\mathbf{x}) = \mathbf{0}$$

$$\mathbf{h}(\mathbf{x}) \leq \mathbf{0}$$

$$\boldsymbol{\nu} \geq \mathbf{0}$$

$$\boldsymbol{\nu}^T \mathbf{h}(\mathbf{x}) = 0$$

at an optimal point, we must have that for either the  $i$ th inequality constraint is active, so  $h_i(\mathbf{x}) = 0$  or it is inactive, but its Lagrange multiplier  $\nu_i = 0$ .



