

CS 450: Numerical Analysis¹

Numerical Optimization

University of Illinois at Urbana-Champaign

¹*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:

Local Minima and Convexity

- ▶ Without knowledge of the analytical form of the function, numerical optimization methods at best achieve convergence to a *local* rather than *global* minimum:

- ▶ A set is *convex* if it includes all points on any line, while a function is (strictly) convex if its (unique) local minimum is always a global minimum:

Optimality Conditions

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

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

- ▶ *Critical points* \boldsymbol{x} satisfy $\nabla f(\boldsymbol{x}) = \mathbf{0}$ 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 consider 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:

Convergence of Steepest Descent

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

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

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 $\mathbf{x}_k - \mathbf{x}_{k-1}$):

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

Krylov Optimization

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

Newton's Method

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

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

x_k - current guess

$$f(x) \approx \tilde{f}(x) = \tilde{f}(x_k + s_k) = f(x_k) + \nabla f(x_k) s_k + \frac{1}{2} s_k^T H_f(x_k) s_k$$

find s_k so that

$$\nabla_{s_k} f(x_k + s_k) = 0$$

solve
↓

$$\nabla f(x_k) + H_f(x_k) s_k = 0 \Rightarrow H_f(x_k) s_k = -\nabla f(x_k)$$

Quasi-Newton Methods

- ▶ **Quasi-Newton** methods compute approximations to the Hessian at each step:

direction s_k can lead to increase if $f(x_k + \underset{\uparrow}{\alpha} s_k)$

$$\underset{\uparrow}{B_k} s_k = -\nabla f(x_k) \quad \left| \quad x_{k+1} = x_k + \underset{\uparrow}{\alpha} s_k \quad \text{for any } \alpha > 0\right.$$

approximate Hessian | line search

- ▶ The **BFGS** method is a secant update method, similar to Broyden's method:

- maintain symmetry of B_k
 - cheap update of B_{k+1} from B_k
 - easy to update B_{k+1}^{-1} from B_k^{-1}
- $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$
 $s_k = x_{k+1} - x_k$

$$\text{BFGS} \Rightarrow B_{k+1} = B_k + \frac{y_k y_k^T}{s_k^T y_k} - \frac{B_k s_k s_k^T B_k}{s_k^T B_k s_k}$$

Nonlinear Least Squares

- ▶ An important special case of multidimensional optimization is nonlinear least squares, the problem of fitting a nonlinear function $f_x(t)$ so that $f_x(t_i) \approx y_i$:

$$f_{[x_1, x_2]}(t) = x_1 \sin(x_2 t) \quad \text{so that} \quad \begin{bmatrix} f_{[x_1, x_2]}(1) \\ \vdots \\ f_{[x_1, x_2]}(t) \end{bmatrix} \approx \begin{bmatrix} 3.2 \\ \vdots \\ 6.4 \end{bmatrix} \} y$$

- ▶ We can cast nonlinear least squares as an optimization problem and solve it by Newton's method: $r_i(x) = f_x(t_i) - y_i$

$$\varphi(x) = \frac{1}{2} \|r(x)\|_2^2 = \frac{1}{2} r(x)^T r(x) \Rightarrow \text{objective function}$$
$$x^* = \underset{x}{\operatorname{arg\,min}} \varphi(x)$$

$$\nabla \varphi(x) = J_r(x) r(x)$$

$$H_{\varphi}(x) = J_r^T(x) J_r(x) + \sum_{i=1} r_i(x) H_{r_i}(x)$$

$$\text{Newton's method} \quad H_{\varphi}(x_k) s_k = -\nabla \varphi(x_k)$$

Gauss-Newton Method

- ▶ The Hessian for nonlinear least squares problems has the form:

$$H_{\mathcal{L}}(x) = J_r(x)^T J_r(x) + \sum_i r_i H_r(x)$$

↑
goes to zero

- ▶ The *Gauss-Newton* method is Newton iteration with an approximate Hessian:

$$H_{\mathcal{L}}(x) \approx J_r(x)^T J_r(x)$$
$$J_r(x_k)^T J_r(x_k) s_k = - \underbrace{J_r(x_k)^T r(x_k)}_{\nabla \mathcal{L}(x_k)} \quad \left| \begin{array}{l} \text{linear least squares} \\ \downarrow \\ J_r(x_k) s_k \approx -r(x_k) \end{array} \right.$$

- ▶ The Levenberg-Marquardt method incorporates Tykhonov regularization into the linear least squares problems within the Gauss-Newton method.

Constrained Optimization Problems

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

$$\min_x f(x) \quad \text{subject to} \quad \underbrace{g(x) = 0}_{\text{equality}} \quad \text{and} \quad \underbrace{h(x) \leq 0}_{\text{inequality}}$$

↑
easier than ↗

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

- ▶ **sequential quadratic programming**: → sequence of QPs that are unconstrained
- ▶ **penalty-based methods**: → sequence of unconstrained but more ill-conditioned problems
- ▶ **active set methods**: → inequality → equality constraints

Sequential Quadratic Programming

- λ_0 \triangleright **Sequential quadratic programming** (SQP) corresponds to using Newton's method to solve the equality constrained optimality conditions, by finding critical points of the Lagrangian function $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})$,

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

- \triangleright At each iteration, SQP computes $\begin{bmatrix} \mathbf{x}_{k+1} \\ \boldsymbol{\lambda}_{k+1} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_k \\ \boldsymbol{\lambda}_k \end{bmatrix} + \begin{bmatrix} \mathbf{s}_k \\ \boldsymbol{\delta}_k \end{bmatrix}$ by solving

$$\underbrace{\mathbf{H}_{\mathcal{L}}(\mathbf{x}_k, \boldsymbol{\lambda}_k)}_{\begin{bmatrix} \mathbf{B}(\mathbf{x}_k, \boldsymbol{\lambda}_k) & \mathbf{J}_g^T(\mathbf{x}_k) \\ \mathbf{J}_g(\mathbf{x}_k) & 0 \end{bmatrix}} \begin{bmatrix} \mathbf{s}_k \\ \boldsymbol{\delta}_k \end{bmatrix} = -\nabla \mathcal{L}(\mathbf{x}_k, \boldsymbol{\lambda}_k)$$

$$\begin{bmatrix} \mathbf{B}(\mathbf{x}_k, \boldsymbol{\lambda}_k) & \mathbf{J}_g^T(\mathbf{x}_k) \\ \mathbf{J}_g(\mathbf{x}_k) & 0 \end{bmatrix}$$

where

$$\mathbf{B}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{H}_f(\mathbf{x}) + \sum_i \lambda_i \mathbf{H}_{g_i}(\mathbf{x})$$

Inequality Constrained Optimality Conditions

► The **Karush-Kuhn-Tucker (KKT)** conditions hold for local minima of a problem with equality and inequality constraints, the key conditions are \sum

- First, any minima x^* must be a feasible point, so $g(x^*) = 0$ and $h(x^*) \geq 0$.
- We say the i th inequality constraint is **active** at a minima x^* if $h_i(x^*) = 0$.
- The collection of equality constraints and active inequality constraints q^* , satisfies $q^*(x^*) = 0$.
- The negative gradient of the objective function at the minima must be in the row span of the Jacobian of this collection of constraints:

$$-\nabla f(x^*) = J_{q^*}^T(x^*)\lambda^* \quad \text{where } \lambda^* \text{ are Lagrange multipliers of constraints in } q^*.$$

► To use SQP for an inequality constrained optimization problem, consider at each iteration an **active set** of constraints:

active set method at step k , we have x_k, λ_k

$$q_k = \begin{bmatrix} g(x) \\ h_2(x) \\ h_4(x) \end{bmatrix} \quad \text{e.g. if } h_2(x) = 0 \text{ and } h_4(x) \geq 0$$

\uparrow exactly satisfied $\underbrace{\hspace{10em}}_{\text{violated}}$

Lagrange multipliers corresp. to q_k

SQP step on $L_k(x, \bar{\lambda}) = f(x) + \sum \lambda^T q_k(x)$

Penalty Functions

- ▶ Alternatively, we can reduce constrained optimization problems to unconstrained ones by modifying the objective function. *Penalty* functions are effective for equality constraints $g(x) = 0$:

$$\mathcal{Q}_p(x) = f(x) + \underset{\substack{\uparrow \\ \text{parameter}}}{p} \langle g(x), g(x) \rangle$$

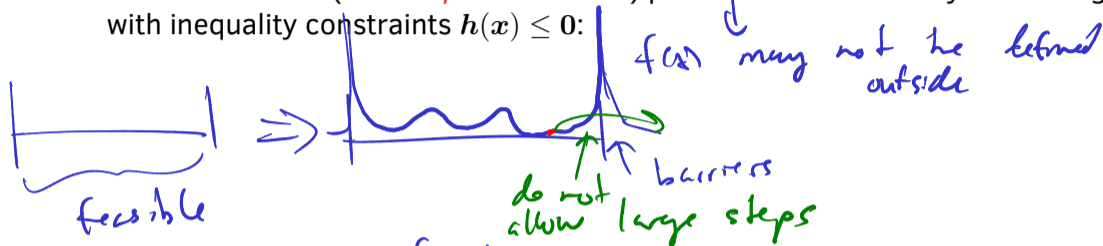
solve sequence of unconstrained problems $\min_x \mathcal{Q}_p(x)$
with $p \rightarrow \infty$, so that $x \rightarrow x^*$

- ▶ The augmented Lagrangian function provides a more numerically robust approach:

$$\mathcal{L}(x, \lambda) = f(x) + \lambda^T g(x) + p \langle g(x), g(x) \rangle$$

Barrier Functions

- **Barrier functions (interior point methods)** provide an effective way of working with inequality constraints $h(x) \leq 0$:



inverse barrier function

$$\ell(x) = f(x) - \mu \frac{1}{h(x)}$$

$$\left. \begin{array}{l} h(x) \leq 0 \\ \text{so as } h(x) \rightarrow 0 \\ \mu \frac{1}{h(x)} \rightarrow -\infty \end{array} \right|$$

take $\mu \rightarrow 0$

so that $\underset{x}{\operatorname{argmin}} (\ell(x)) \rightarrow \underset{x}{\operatorname{argmin}} (f(x)) \text{ s.t. } h(x) \leq 0$

logarithmic barrier function: $\ell(x) = f(x) - \mu \log(-h(x))$