

CS 450: Numerical Analysis

Lecture 17

Chapter 6 Numerical Optimization

Constrained Optimization and Quadratic Programming

Edgar Solomonik

Department of Computer Science
University of Illinois at Urbana-Champaign

March 14, 2018

Constrained Optimization Problems

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

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

*When f is quadratic, while h, g is linear, this is a *quadratic optimization problem*.*

- ▶ Generally, we will seek to reduce constrained optimization problems to a series of unconstrained optimization problems:
 - ▶ *sequential quadratic programming*: solve an unconstrained quadratic program at each iteration,
 - ▶ *penalty-based methods*: solve a series of more complicated (more ill-conditioned) unconstrained optimization problems,
 - ▶ *active set methods*: define sequence of optimization problems with inequality constrained ignored or treated as equality constraints.

Lagrangian Duality

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

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

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}_1^T \mathbf{h}(\mathbf{x}) + \boldsymbol{\lambda}_2^T \mathbf{g}(\mathbf{x})$$

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

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

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

$$\max \left(\boldsymbol{\lambda}^*, \begin{bmatrix} \mathbf{h}(\mathbf{x}) \\ \mathbf{g}(\mathbf{x}) \end{bmatrix} \right) = \mathbf{0}$$

when $\lambda_i^ = 0$, we say the i th constraint is inactive.*

Constrained Optimality

- ▶ In equality-constrained optimization $\mathbf{g}(\mathbf{x}) = \mathbf{0}$, minimizers \mathbf{x}^* are on the border of the feasible region (set of points satisfying constraints), in which case we must ensure any direction of decrease of f from \mathbf{x}^* leads to an infeasible point, which gives us the condition:

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

$\boldsymbol{\lambda}$ are referred to as the Lagrange multipliers.

- ▶ Seek critical points in the Lagrangian function $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})$, described by the nonlinear equation,

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

Seeking $\boldsymbol{\lambda}$ that maximizes the global minimum of $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ defines the *dual optimization problem*.

Sequential Quadratic Programming

- ▶ *Sequential quadratic programming (SQP)* corresponds to using Newton's method to solve the nonlinear equations,

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

At each iteration, we compute $\begin{bmatrix} \mathbf{x}_{k+1} \\ \boldsymbol{\lambda}_{k+1} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_k \\ \boldsymbol{\lambda}_k \end{bmatrix} + \mathbf{s}_k$ by solving

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

where

$$\mathbf{J}_{\mathcal{L}}(\mathbf{x}_k, \boldsymbol{\lambda}_k) = \begin{bmatrix} \mathbf{B}(\mathbf{x}_k, \boldsymbol{\lambda}_k) & \mathbf{J}_g^T(\mathbf{x}_k) \\ \mathbf{J}_g(\mathbf{x}_k) & \mathbf{0} \end{bmatrix} \quad \text{with} \quad \mathbf{B}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{H}_f(\mathbf{x}) + \sum_{i=1}^m \lambda_i \mathbf{H}_{g_i}(\mathbf{x})$$

Quadratic Programming Problems

- ▶ An equality-constrained quadratic programming problem has the form

$$\min_{\mathbf{x}} f(\mathbf{x}), \quad f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{c}^T \mathbf{x} \quad \text{subject to} \quad \mathbf{A} \mathbf{x} = \mathbf{b}$$

Its first-order optimality condition in the unconstrained case is

$$\mathbf{0} = \nabla f(\mathbf{x}) = \mathbf{Q} \mathbf{x} + \mathbf{c} \quad \Rightarrow \quad \mathbf{Q} \mathbf{x} = -\mathbf{c}$$

and in the constrained case, becomes

$$\begin{bmatrix} \mathbf{Q} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\lambda} \end{bmatrix} = - \begin{bmatrix} \mathbf{c} \\ \mathbf{b} \end{bmatrix}$$

This allows us to observe, that at the k th iteration SQP solves a QP with $\mathbf{Q} = \mathbf{B}(\mathbf{x}_k, \boldsymbol{\lambda}_k)$, $\mathbf{A} = \mathbf{J}_g(\mathbf{x}_k)$, $\mathbf{c} = \nabla f(\mathbf{x}_k) + \mathbf{J}_g^T(\mathbf{x}_k) \boldsymbol{\lambda}_k$, and $\mathbf{b} = \mathbf{g}(\mathbf{x}_k)$.

Steepest Descent for Quadratic Programming

- ▶ Near the minima, all smooth nonlinear programming problems look like quadratic programming problems, where \mathbf{Q} converges to the Hessian at the minima, $\mathbf{H}_f(\mathbf{x}^*)$.
- ▶ Consequently, we can analyze local convergence of methods by considering their convergence for a QP, e.g. for steepest descent:

Assume $\mathbf{b} = \mathbf{0}$ and $\mathbf{c} = \mathbf{0}$, then we have $\nabla f(\mathbf{x}) = \mathbf{Q}\mathbf{x}$. Consequently, steepest descent will seek a step-size α_k to perform the step

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{Q}\mathbf{x}_k = (\mathbf{I} - \alpha_k \mathbf{Q})\mathbf{x}_k$$

for this fixed-point iteration to converge, it suffices that $\alpha_k < 2/\lambda_{\max}(\mathbf{Q})$. The optimal value can be shown to be $\alpha^ = 2/(\lambda_{\max}(\mathbf{Q}) + \lambda_{\min}(\mathbf{Q}))$, leading to convergence rate*

$$\|e_{k+1}\| = \frac{\kappa(\mathbf{Q}) - 1}{\kappa(\mathbf{Q}) + 1} \|e_k\|$$

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 (using *momentum*):

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1})$$

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

$$\|\mathbf{e}_{k+1}\| = \frac{\sqrt{\kappa(\mathbf{Q})} - 1}{\sqrt{\kappa(\mathbf{Q})} + 1} \|\mathbf{e}_k\|$$

Nesterov's gradient optimization method is another instance of an extrapolation method, and provides further optimality guarantees.

Conjugate Gradient Method

- ▶ The *conjugate gradient method* is capable of making the optimal choice of α_k and β_k at each iteration of an extrapolation method:

$$(\alpha_k, \beta_k) = \underset{\alpha_k, \beta_k}{\operatorname{argmin}} \left[f\left(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1})\right) \right]$$

For quadratic programming problems, conjugate gradient is an optimal 1st order method, converging in n iterations.

- ▶ Generally conjugate gradient methods perform a sequence of line minimizations in n directions that are Q -orthogonal:

*A **parallel tangents** implementation of the method proceeds as follows*

- ▶ *Perform a step of steepest descent to generate $\hat{\mathbf{x}}_k$ from \mathbf{x}_k*
- ▶ *Generate \mathbf{x}_{k+1} by minimizing over the line passing through \mathbf{x}_k and $\hat{\mathbf{x}}_k$*

Each conjugate direction is guaranteed to be Q -orthogonal to previous directions by a recurrence argument similar to that of Lanczos iteration.