

# CS 450: Numerical Analysis

Lecture 18

Chapter 6 Numerical Optimization

Conjugate Gradient and Constrained Optimization

Edgar Solomonik

Department of Computer Science  
University of Illinois at Urbana-Champaign

March 27, 2018

## Sequential Quadratic Programming

- ▶ *Sequential quadratic programming (SQP)* solves a constrained quadratic program at the  $k$ th step:

*We approximate Lagrangian function  $\mathcal{L}_f(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})$  by taking its Taylor expansion in  $\mathbf{x}$  and  $\boldsymbol{\lambda}$ :*

$$\begin{aligned} q(\mathbf{x}_k + \mathbf{s}, \boldsymbol{\lambda}_k + \boldsymbol{\delta}) = & \mathcal{L}_f(\mathbf{x}_k, \boldsymbol{\lambda}_k) + \mathbf{s}^T (\nabla f(\mathbf{x}_k) + \mathbf{J}_g^T(\mathbf{x}_k) \boldsymbol{\lambda}_k) + \frac{1}{2} \mathbf{s}^T \mathbf{B}(\mathbf{x}_k) \mathbf{s} \\ & + \boldsymbol{\delta} (\mathbf{J}_g(\mathbf{x}_k) \mathbf{s} + \mathbf{g}(\mathbf{x}_k)) \end{aligned}$$

*where we can minimize over  $\mathbf{s}$  and  $\boldsymbol{\delta}$  while ignoring the constant term  $\mathcal{L}_f(\mathbf{x}_k, \boldsymbol{\lambda}_k)$ . This unconstrained quadratic program corresponds to the Lagrangian form of the constrained quadratic program*

$$\max_{\mathbf{s}} \mathbf{s}^T (\nabla f(\mathbf{x}_k) + \mathbf{J}_g^T(\mathbf{x}_k) \boldsymbol{\lambda}_k) + \frac{1}{2} \mathbf{s}^T \mathbf{B}(\mathbf{x}_k) \mathbf{s}$$

*with constraint  $\mathbf{J}_g(\mathbf{x}_k) \mathbf{s} = -\mathbf{g}(\mathbf{x}_k)$ .*

## Solving Quadratic Programs

- ▶ Newton's method for optimization can solve the quadratic program with a single step:

*Quadratic approximation is exact for an unconstrained QP. However, Newton's method forces us to form and invert  $\mathbf{H}_f(\mathbf{x}_k)$ , which require  $O(n^3)$  cost and  $O(n^2)$  storage.*

- ▶ The *conjugate gradient method* provides an effective way of solving QPs iteratively:

*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  $\mathbf{x}_{k-1}$  and  $\hat{\mathbf{x}}_k$*

*Each conjugate direction is guaranteed to be  $\mathbf{A}$ -orthogonal to previous directions by a recurrence argument similar to that of Lanczos iteration.*

## Conjugate Gradient as a Krylov Subspace Method

- ▶ Generally, Krylov subspaces describe the information available from  $k$  matrix-vector products, and can be used to find an approximation  $\mathbf{x}_k$  to the minima of  $\mathbf{c}^T \mathbf{x} - \mathbf{x}^T \mathbf{A} \mathbf{x}$ :

*Recall the Krylov subspace is defined as  $\mathcal{K}_r = \text{span}(\mathbf{x}_0, \mathbf{A}\mathbf{x}_0, \dots, \mathbf{A}^{r-1}\mathbf{x}_0)$ , so our goal is to select a good  $\mathbf{x}_k \in \mathcal{K}_r$ :*

- ▶ *there is insufficient information to minimize  $\|\mathbf{x}_k - \mathbf{x}^*\|_2^2$*
- ▶ *minimum-residual methods (MINRES/GMRES) minimize residual  $\|\mathbf{c}^T - \mathbf{A}\mathbf{x}_k\|_2^2$*
- ▶ *conjugate gradient methods minimize residual in  $\mathbf{A}^{-1}$  norm*  
$$\|\mathbf{c}^T - \mathbf{A}\mathbf{x}_k\|_{\mathbf{A}^{-1}}^2 = (\mathbf{c}^T - \mathbf{A}\mathbf{x}_k)^T \mathbf{A}^{-1} (\mathbf{c}^T - \mathbf{A}\mathbf{x}_k)$$
- ▶ Conjugate gradient can be derived from vectors generated by the Lanczos algorithm for symmetric (positive-definite)  $\mathbf{A}$ , yielding

$$\mathbf{x}_k = \mathbf{Q}_k \mathbf{T}_k^{-1} \mathbf{e}_1 \|\mathbf{c}\|_2 \quad \text{if} \quad \mathbf{x}_0 = \mathbf{c}$$

*Minimization follows from the fact that we minimize the projection of the residual onto  $\mathbf{Q}$ , namely  $\mathbf{Q}^T \mathbf{r}$ , since*

$$\mathbf{Q}^T (\mathbf{c} - \mathbf{A}\mathbf{x}) = \mathbf{e}_1 \|\mathbf{c}\|_2 - \underbrace{\mathbf{Q}^T \mathbf{A} \mathbf{Q}}_{\mathbf{T}} \mathbf{T}^{-1} \mathbf{e}_1 \|\mathbf{c}\|_2 = \mathbf{0}$$

## Conjugate Gradient Properties

- ▶ Each iteration of conjugate gradient has cost proportional to a matrix-vector product:
  - ▶ *Lanczos generates each vector  $\mathbf{q}_{k+1}$  by orthogonalizing  $\mathbf{A}\mathbf{q}_k$  to  $\mathbf{q}_{k-1}$  and  $\mathbf{q}_k$  (three term recurrence)*
  - ▶ *Note that the directions of change  $\mathbf{q}_{k+1}$  and  $\mathbf{q}_{k-1}$  are  $\mathbf{A}$ -orthogonal since*

$$\mathbf{q}_{k+1}^T \mathbf{A} \mathbf{q}_{k-1} = \mathbf{q}_{k+1}^T (\alpha \mathbf{q}_k - \beta \mathbf{q}_{k-1} - \gamma \mathbf{q}_{k-2}) = 0$$

- ▶ *Conjugate gradient method obtained by transforming recurrence for  $\mathbf{q}_k$  into recurrence for  $\mathbf{x}_k$ , requiring storage of 4 vectors*
- ▶ Conjugate gradient is especially efficient when the matrix has a sparse or implicit representation:  
*Requires only ability to perform matrix-vector product and store  $O(1)$  vectors.*

## Active Set Methods

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

*A constraint is active if the current iterate is at a boundary satisfies the inequality constraint as an equality (i.e. current position is at boundary generated by constraint).*

- ▶ 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  $\nu_i = 0$ .

## Penalty Functions

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

$$\phi_\rho(\mathbf{x}) = f(\mathbf{x}) + \frac{1}{2}\rho\mathbf{g}(\mathbf{x})^T\mathbf{g}(\mathbf{x})$$

*is a simple merit function, and its solutions  $\mathbf{x}_\rho^*$  satisfy  $\lim_{\rho \rightarrow \infty} \mathbf{x}_\rho^* = \mathbf{x}^*$ . However, the Hessian of  $\phi_\rho$  becomes increasingly ill-conditioned for large  $\rho$ , leading to slow convergence.*

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

$$\mathcal{L}_\rho(\mathbf{x}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T\mathbf{g}(\mathbf{x}) + \frac{1}{2}\rho\mathbf{g}(\mathbf{x})^T\mathbf{g}(\mathbf{x})$$

## Barrier Functions

- ▶ A drawback of penalty function methods is that they can produce infeasible approximate solutions, which is problematic if the objective function is only defined in the feasible region:

*Moreover, the merit and augmented Lagrangian functions only make sense for equality constraints.*

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

*Inverse barrier function:*

$$\phi_{\mu}(\mathbf{x}) = f(\mathbf{x}) - \mu \sum_{i=1}^m \frac{1}{h_i(\mathbf{x})}$$

*Logarithmic barrier function:*

$$\phi_{\mu}(\mathbf{x}) = f(\mathbf{x}) - \mu \sum_{i=1}^m \log(-h_i(\mathbf{x}))$$

*in theory with sufficiently small steps we have  $\mathbf{x}_{\mu}^* \rightarrow \mathbf{x}^*$  as  $\mu \rightarrow 0$*