

STAT 309: MATHEMATICAL COMPUTATIONS I  
FALL 2015  
LECTURE 11

1. LEAST SQUARES WITH LINEAR CONSTRAINTS

- suppose that we wish to fit data as in the least squares problem, except that we are using different functions to fit the data on different subintervals
- a common example is the process of fitting data using cubic splines, with a different cubic polynomial approximating data on each subinterval
- typically it is desired that the functions assigned to each piece form a function that is continuous on the entire interval within which the data lies
- this requires that *constraints* be imposed on the functions themselves
- it is also not uncommon to require that the function assembled from these pieces also has a continuous first or even second derivative, resulting in additional constraints
- the result is a *least squares problem with linear constraints*, as the constraints are applied to coefficients of predetermined functions chosen as a basis for some function space, such as the space of polynomials of a given degree
- the general form of a least squares problem with linear constraints is as follows: we wish to find an  $\mathbf{x} \in \mathbb{R}^n$  that minimizes  $\|\mathbf{Ax} - \mathbf{b}\|_2$ , subject to the constraint  $C^T \mathbf{x} = \mathbf{d}$ , where  $A \in \mathbb{R}^{m \times n}$ ,  $C \in \mathbb{R}^{n \times p}$ ,  $\mathbf{b} \in \mathbb{R}^m$ , and  $\mathbf{d} \in \mathbb{R}^p$  are given

$$\begin{array}{ll} \text{minimize} & \|\mathbf{b} - \mathbf{Ax}\|_2^2 \\ \text{subject to} & C^T \mathbf{x} = \mathbf{d} \end{array} \quad (1.1)$$

- again we will describe three methods, mathematically equivalent but with different numerical properties
- this problem is usually solved using *Lagrange multipliers*, define

$$L(\mathbf{x}, \boldsymbol{\lambda}) = \|\mathbf{b} - \mathbf{Ax}\|_2^2 + 2\boldsymbol{\lambda}^T(C^T \mathbf{x} - \mathbf{d})$$

- in optimization parlance, the function  $L$  is called the *Lagrangian* and  $\boldsymbol{\lambda}^T = [\lambda_1, \dots, \lambda_p]^T$  is the vector of Lagrange multipliers
- setting derivative with respect to  $\mathbf{x}$  to zero yields

$$\mathbf{0} = \nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}) = 2(A^T \mathbf{Ax} - A^T \mathbf{b} + C\boldsymbol{\lambda})$$

and so

$$A^T \mathbf{Ax} + C\boldsymbol{\lambda} = A^T \mathbf{b} \quad (1.2)$$

- note that  $\mathbf{0} = \nabla_{\boldsymbol{\lambda}} L(\mathbf{x}, \boldsymbol{\lambda})$  just gives us back the constraint

$$C^T \mathbf{x} = \mathbf{d} \quad (1.3)$$

- in optimization parlance, (1.2) and (1.3) are collectively called the *KKT conditions*
- method 1: together (1.2) and (1.3) give the system

$$\begin{bmatrix} A^T A & C \\ C^T & 0 \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\lambda} \end{bmatrix} = \begin{bmatrix} A^T \mathbf{b} \\ \mathbf{d} \end{bmatrix}$$

- solving this linear system gives us a solution to (1.1)

- this method preserves the sparsity of  $C$  but involves a coefficient matrix of size  $(n + p) \times (n + p)$ , larger than the next two methods
- method 2: if  $A$  has full column rank, then from  $A^T A \mathbf{x} = A^T \mathbf{b} - C \boldsymbol{\lambda}$ , we see that we can first compute  $\mathbf{x} = \hat{\mathbf{x}} - (A^T A)^{-1} C \boldsymbol{\lambda}$  where  $\hat{\mathbf{x}}$  is the solution to the unconstrained least squares problem

$$\hat{\mathbf{x}} \in \operatorname{argmin} \|A \mathbf{x} - \mathbf{b}\|_2$$

- then from the equation  $C^T \mathbf{x} = \mathbf{d}$  we obtain the  $p \times p$  linear system

$$C^T (A^T A)^{-1} C \boldsymbol{\lambda} = C^T \hat{\mathbf{x}} - \mathbf{d} \quad (1.4)$$

which we can then solve for  $\boldsymbol{\lambda}$

- this works because  $A^T A \hat{\mathbf{x}} = A^T \mathbf{b}$  and therefore

$$A^T A \mathbf{x} = A^T \mathbf{b} - C \boldsymbol{\lambda}$$

- the actual algorithm uses two QR factorization and does not actually require solving a system involving  $A^T A$ 
  - compute full-rank QR

$$A = Q \begin{bmatrix} R \\ 0 \end{bmatrix}$$

with nonsingular  $R \in \mathbb{R}^{n \times n}$

- solve the unconstrained least squares problem  $\min \|A \mathbf{x} - \mathbf{b}\|_2$  for  $\hat{\mathbf{x}}$
- form  $W = R^{-T} C$
- compute  $W = Q_1 R_1$ , the QR factorization of  $W$
- solve  $R_1^T R_1 \boldsymbol{\lambda} = \boldsymbol{\eta} = C^T \hat{\mathbf{x}} - \mathbf{d}$  for  $\boldsymbol{\lambda}$
- set  $\mathbf{x} = \hat{\mathbf{x}} - (A^T A)^{-1} C \boldsymbol{\lambda}$
- this works because

$$\begin{aligned} R_1^T R_1 &= (Q_1^T W)^T (Q_1^T W) \\ &= W^T Q_1 Q_1^T W \\ &= C^T R^{-1} R^{-T} C \\ &= C^T (R^T R)^{-1} C \\ &= C^T (R^T Q^T Q R)^{-1} C \\ &= C^T (A^T A)^{-1} C \end{aligned}$$

- method 2 is not the most practical since it has more unknowns than the unconstrained least squares problem, which is odd because the constraints should have the effect of eliminating unknowns, not adding them
- method 3: suppose  $p \leq n$ , then computing the QR factorization of  $C$  yields

$$C = Q_2 \begin{bmatrix} R_2 \\ 0 \end{bmatrix}$$

where  $R_2$  is a  $p \times p$  upper triangular matrix

- then the constraint  $C^T \mathbf{x} = \mathbf{d}$  takes the form

$$R_2^T \mathbf{u} = \mathbf{d}, \quad Q_2^T \mathbf{x} = \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}$$

- then

$$\begin{aligned}
\|\mathbf{b} - A\mathbf{x}\|_2 &= \|\mathbf{b} - AQ_2Q_2^\top \mathbf{x}\|_2 \\
&= \left\| \mathbf{b} - \tilde{A} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} \right\|_2, \quad \tilde{A} = AQ_2 \\
&= \left\| \mathbf{b} - [\tilde{A}_1 \quad \tilde{A}_2] \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} \right\|_2 \\
&= \|\mathbf{b} - \tilde{A}_1\mathbf{u} - \tilde{A}_2\mathbf{v}\|_2
\end{aligned}$$

- thus we can obtain  $\mathbf{x}$  by the following algorithm:
  - compute the QR factorization of  $C$
  - compute  $\tilde{A} = AQ_2$
  - solve  $R_2^\top \mathbf{u} = \mathbf{d}$
  - solve the new least squares problem of minimizing  $\|(\mathbf{b} - \tilde{A}_1\mathbf{u}) - \tilde{A}_2\mathbf{v}\|_2$
  - compute

$$\mathbf{x} = Q_2 \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}$$

- method 3 has the advantage that there are fewer unknowns in each system that needs to be solved, and also that  $\kappa_2(\tilde{A}_2) \leq \kappa_2(A)$
- the drawback is that sparsity can be destroyed

## 2. COMPUTING THE QR FACTORIZATION

- there are two common ways to compute the QR decomposition:
  - using *Householder matrices*, developed by Alston S. Householder
  - using *Givens rotations*, also known as *Jacobi rotations*, used by Wallace Givens and originally invented by Jacobi for use with in solving the symmetric eigenvalue problem in 1846
  - the Gram–Schmidt or modified Gram–Schmidt orthogonalization discussed in previous lecture works in principle but has numerical stability issues and are not usually used
- roughly speaking, Gram–Schmidt applies a sequence of triangular matrices to orthogonalize  $A$  (i.e., transform  $A$  into an orthogonal matrix  $Q$ ),

$$AR_1^{-1}R_2^{-1} \cdots R_{n-1}^{-1} = Q$$

whereas Householder and Givens QR apply a sequence of orthogonal matrices to triangularize  $A$  (i.e., transform  $A$  into an upper triangular matrix  $R$ ),

$$Q_{n-1}^\top \cdots Q_2^\top Q_1^\top A = R$$

- orthogonal transformations are highly desirable in algorithms as they preserve lengths and therefore do not blow up the errors present at every stage of the computation

## 3. ORTHOGONALIZATION USING GIVENS ROTATIONS

- we illustrate the process in the case where  $A$  is a  $2 \times 2$  matrix
- in Gaussian elimination, we compute  $L^{-1}A = U$  where  $L^{-1}$  is unit lower triangular and  $U$  is upper triangular, specifically,

$$\begin{bmatrix} 1 & 0 \\ m_{21} & 1 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} = \begin{bmatrix} a_{11}^{(2)} & a_{12}^{(2)} \\ 0 & a_{22}^{(2)} \end{bmatrix}, \quad m_{21} = -\frac{a_{21}}{a_{11}}$$

- by contrast, the QR decomposition takes the form

$$\begin{bmatrix} \gamma & \sigma \\ -\sigma & \gamma \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} = \begin{bmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{bmatrix}$$

where  $\gamma^2 + \sigma^2 = 1$

- from the relationship  $-\sigma a_{11} + \gamma a_{21} = 0$  we obtain

$$\begin{aligned} \gamma a_{21} &= \sigma a_{11} \\ \gamma^2 a_{21}^2 &= \sigma^2 a_{11}^2 = (1 - \gamma^2) a_{11}^2 \end{aligned}$$

which yields

$$\gamma = \pm \frac{a_{11}}{\sqrt{a_{21}^2 + a_{11}^2}}$$

- it is conventional to choose the + sign
- then, we obtain

$$\sigma^2 = 1 - \gamma^2 = 1 - \frac{a_{11}^2}{a_{21}^2 + a_{11}^2} = \frac{a_{21}^2}{a_{21}^2 + a_{11}^2},$$

or

$$\sigma = \pm \frac{a_{21}}{\sqrt{a_{21}^2 + a_{11}^2}}$$

- again, we choose the + sign
- as a result, we have

$$r_{11} = a_{11} \frac{a_{11}}{\sqrt{a_{21}^2 + a_{11}^2}} + a_{21} \frac{a_{21}}{\sqrt{a_{21}^2 + a_{11}^2}} = \sqrt{a_{21}^2 + a_{11}^2}$$

- the matrix

$$Q^T = \begin{bmatrix} \gamma & \sigma \\ -\sigma & \gamma \end{bmatrix}$$

is called a *Givens rotation*

- it is called a rotation because it is orthogonal, and therefore length-preserving, and also because there is an angle  $\theta$  such that  $\sin \theta = \sigma$  and  $\cos \theta = \gamma$ , and its effect is to rotate a vector through the angle  $\theta$
- in particular,

$$\begin{bmatrix} \gamma & \sigma \\ -\sigma & \gamma \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} \rho \\ 0 \end{bmatrix}$$

where  $\rho = \sqrt{\alpha^2 + \beta^2}$ ,  $\alpha = \rho \cos \theta$  and  $\beta = \rho \sin \theta$

- it is easy to verify that the product of two rotations is itself a rotation
- now, in the case where  $A$  is an  $n \times n$  matrix, suppose that we are given the vector

$$\begin{bmatrix} \times & \cdots & \times & \alpha & \times & \cdots & \times & \beta & \times & \cdots & \times \end{bmatrix}^T \in \mathbb{R}^n,$$

then

$$\begin{bmatrix}
 1 & & & & & & & & & & \\
 & \ddots & & & & & & & & & \\
 & & 1 & & & & & & & & \\
 & & & \gamma & & & & & & & \\
 & & & & 1 & & & \sigma & & & \\
 & & & & & \ddots & & & & & \\
 & & & & & & 1 & & & & \\
 & & & -\sigma & & & & \gamma & & & \\
 & & & & & & & & 1 & & \\
 & & & & & & & & & \ddots & \\
 & & & & & & & & & & 1
 \end{bmatrix}
 \begin{bmatrix}
 \times \\
 \vdots \\
 \times \\
 \alpha \\
 \times \\
 \vdots \\
 \times \\
 \beta \\
 \times \\
 \vdots \\
 \times
 \end{bmatrix}
 =
 \begin{bmatrix}
 \times \\
 \vdots \\
 \times \\
 \rho \\
 \times \\
 \vdots \\
 \times \\
 0 \\
 \times \\
 \vdots \\
 \times
 \end{bmatrix}$$

- so, in order to transform  $A$  into an upper triangular matrix  $R$ , we can find a product of rotations  $Q$  such that  $Q^T A = R$
- it is easy to see that  $O(n^2)$  rotations are required