

Mathematical Tripos Part IB: Lent 2020

Numerical Analysis – Lecture 16¹

Householder reflections Let $\mathbf{u} \in \mathbb{R}^m \setminus \{\mathbf{0}\}$. The $m \times m$ matrix $I - 2\frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2}$ is called a *Householder reflection*. Each such matrix is symmetric and orthogonal, since

$$\left(I - 2\frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2}\right)^\top \left(I - 2\frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2}\right) = \left(I - 2\frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2}\right)^2 = I - 4\frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2} + 4\frac{\mathbf{u}(\mathbf{u}^\top\mathbf{u})\mathbf{u}^\top}{\|\mathbf{u}\|^4} = I.$$

Householder reflections offer an alternative to Given rotations in the calculation of a QR factorization.

Householder algorithm Our goal is to multiply an $m \times n$ matrix A by a sequence of Householder reflections so that each product induces zeros under the diagonal in an entire column.

At the first step we seek a reflection that transforms the first column \mathbf{a}_1 of A to a multiple of \mathbf{e}_1 . Since the Householder reflection is orthogonal (it preserves Euclidean norm) the latter has to be $\pm\|\mathbf{a}_1\|\mathbf{e}_1$ where we are free to choose the sign. The Householder reflection that does this operation is given by the choice of vector $\mathbf{u} = \mathbf{a}_1 - (\pm\|\mathbf{a}_1\|\mathbf{e}_1)$. For numerical stability the sign is usually chosen to be $-\text{sign}(A_{11})$.

More generally, at the beginning of the k 'th step of the algorithm, the columns 1 to $k-1$ have been processed and have zeros under their diagonal element. Our goal is to find a Householder reflection that will induce zeros under the diagonal element of the k 'th column. To do so we use a block orthogonal matrix $\begin{bmatrix} I & \mathbf{0} \\ \mathbf{0} & H \end{bmatrix}$ where I is a $(k-1) \times (k-1)$ identity matrix, and H is a $(m-k+1) \times (m-k+1)$ Householder reflection associated with the choice $\tilde{\mathbf{u}} = \tilde{\mathbf{a}}_k + \text{sign}(A_{kk})\|\tilde{\mathbf{a}}_k\|\tilde{\mathbf{e}}_1$, where $\tilde{\mathbf{a}}_k$ is the vector of size $m-k+1$ consisting of the entries of A under the diagonal in the k 'th column, and $\tilde{\mathbf{e}}_1$ is the vector of size $m-k+1$ with a 1 in the first position and zero elsewhere.

To summarize it is convenient to use the (Matlab-style) notation where $A_{k:m,j}$ indicates the vector of size $m-k+1$ obtained from rows k, \dots, m of column j of A . Then the algorithm can be written as follows:

Given $A \in \mathbb{R}^{m \times n}$ with $m \geq n$. For $k = 1$ to n :

- Let $\tilde{\mathbf{a}}_k = A_{k:m,k} \in \mathbb{R}^{m-k+1}$
- Let $\tilde{\mathbf{e}}_1$ be the vector of size $m-k+1$ with a 1 in the first position and zero elsewhere.
- Let $\tilde{\mathbf{u}} = \tilde{\mathbf{a}}_k + \text{sign}(A_{kk})\|\tilde{\mathbf{a}}_k\|\tilde{\mathbf{e}}_1$
- For each column $j = k, \dots, n$ update $A_{k:m,j} = A_{k:m,j} - 2(\tilde{\mathbf{u}}^\top A_{k:m,j})\tilde{\mathbf{u}}/\|\tilde{\mathbf{u}}\|^2$.

Example ($k = 3$, assuming the first two columns have already been processed)

$$A = \begin{bmatrix} 2 & 4 & 7 \\ 0 & 3 & -1 \\ 0 & 0 & 2 \\ 0 & 0 & 1 \\ 0 & 0 & -2 \end{bmatrix} \rightarrow \tilde{\mathbf{a}}_3 = \begin{bmatrix} 2 \\ 1 \\ -2 \end{bmatrix}, \quad \tilde{\mathbf{u}} = \begin{bmatrix} 5 \\ 1 \\ -2 \end{bmatrix} \rightarrow \begin{bmatrix} 2 & 4 & 7 \\ 0 & 3 & -1 \\ 0 & 0 & -3 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

Calculation of Q Like for the case of Givens algorithm, the matrix Q is not explicitly formed. To form Q explicitly we start with $\Omega = I$ initially and, for each step we replace Ω , by $\left(I - 2\frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2}\right)\Omega = \Omega - \frac{2}{\|\mathbf{u}\|^2}\mathbf{u}(\mathbf{u}^\top\Omega)$ where $\mathbf{u} = \begin{bmatrix} \mathbf{0} \\ \tilde{\mathbf{u}} \end{bmatrix}$ is obtained from $\tilde{\mathbf{u}}$ by adding $k-1$ zeros above it². However, if we require just the vector $\mathbf{c} = Q^\top \mathbf{b}$, say, rather than the matrix Q , then we set initially $\mathbf{c} = \mathbf{b}$ and in each stage replace \mathbf{c} by $\left(I - 2\frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2}\right)\mathbf{c} = \mathbf{c} - 2\frac{\mathbf{u}^\top \mathbf{c}}{\|\mathbf{u}\|^2}\mathbf{u}$.

¹Corrections and suggestions to these notes should be emailed to h.fawzi@damp.cam.ac.uk.

²Indeed, note that the reflection $I - 2\mathbf{u}\mathbf{u}^\top/\|\mathbf{u}\|^2$ is the same as the block orthogonal matrix $\begin{bmatrix} I & \mathbf{0} \\ \mathbf{0} & H \end{bmatrix}$ where H is the Householder reflection corresponding to $\tilde{\mathbf{u}}$.

Givens or Householder? If A is dense, it is in general more convenient to use Householder reflections. Givens rotations come into their own, however, when A has many leading zeros in its rows. E.g., if an $n \times n$ matrix A consists of zeros underneath the first subdiagonal, they can be ‘rotated away’ in just $n - 1$ Givens rotations, at the cost of $\mathcal{O}(n^2)$ operations!

5.3 Linear least squares

Statement of the problem Suppose that an $m \times n$ matrix A and a vector $\mathbf{b} \in \mathbb{R}^m$ are given. The equation $A\mathbf{x} = \mathbf{b}$, where $\mathbf{x} \in \mathbb{R}^n$ is unknown, has in general no solution (if $m > n$) or an infinity of solutions (if $m < n$). Problems of this form occur frequently when we collect m observations (which, typically, are prone to measurement error) and wish to exploit them to form an n -variable linear model, where $n \ll m$. (In statistics, this is known as *linear regression*.) Bearing in mind the likely presence of errors in A and \mathbf{b} , we seek $\mathbf{x} \in \mathbb{R}^n$ that minimises the Euclidean length $\|A\mathbf{x} - \mathbf{b}\|$. This is the *least squares problem*.

Theorem $\mathbf{x} \in \mathbb{R}^n$ is a solution of the least squares problem iff $A^\top(A\mathbf{x} - \mathbf{b}) = \mathbf{0}$.

Proof. If \mathbf{x} is a solution then it minimises

$$f(\mathbf{x}) := \|A\mathbf{x} - \mathbf{b}\|^2 = \langle A\mathbf{x} - \mathbf{b}, A\mathbf{x} - \mathbf{b} \rangle = \mathbf{x}^\top A^\top A\mathbf{x} - 2\mathbf{x}^\top A^\top \mathbf{b} + \mathbf{b}^\top \mathbf{b}.$$

Hence $\nabla f(\mathbf{x}) = \mathbf{0}$. But $\frac{1}{2}\nabla f(\mathbf{x}) = A^\top A\mathbf{x} - A^\top \mathbf{b}$, hence $A^\top(A\mathbf{x} - \mathbf{b}) = \mathbf{0}$.

Conversely, suppose that $A^\top(A\mathbf{x} - \mathbf{b}) = \mathbf{0}$ and let $\mathbf{u} \in \mathbb{R}^n$. Hence, letting $\mathbf{y} = \mathbf{u} - \mathbf{x}$,

$$\begin{aligned} \|A\mathbf{u} - \mathbf{b}\|^2 &= \langle A\mathbf{x} + A\mathbf{y} - \mathbf{b}, A\mathbf{x} + A\mathbf{y} - \mathbf{b} \rangle = \langle A\mathbf{x} - \mathbf{b}, A\mathbf{x} - \mathbf{b} \rangle + 2\mathbf{y}^\top A^\top(A\mathbf{x} - \mathbf{b}) \\ &\quad + \langle A\mathbf{y}, A\mathbf{y} \rangle = \|A\mathbf{x} - \mathbf{b}\|^2 + \|A\mathbf{y}\|^2 \geq \|A\mathbf{x} - \mathbf{b}\|^2 \end{aligned}$$

and \mathbf{x} is indeed optimal. □

Corollary Optimality of $\mathbf{x} \Leftrightarrow$ the vector $A\mathbf{x} - \mathbf{b}$ is orthogonal to all columns of A .

Normal equations One way of finding optimal \mathbf{x} is by solving the $n \times n$ linear system $A^\top A\mathbf{x} = A^\top \mathbf{b}$; this is the method of *normal equations*. This approach is popular in many applications. However, there are three disadvantages. Firstly, $A^\top A$ might be singular, secondly sparse A might be replaced by a dense $A^\top A$ and, finally, forming $A^\top A$ might lead to loss of accuracy. Thus, suppose that our computer works in the IEEE arithmetic standard (≈ 15 significant digits) and let

$$A = \begin{bmatrix} 10^8 & -10^8 \\ 1 & 1 \end{bmatrix} \quad \Rightarrow \quad A^\top A = \begin{bmatrix} 10^{16} + 1 & -10^{16} + 1 \\ -10^{16} + 1 & 10^{16} + 1 \end{bmatrix} \approx 10^{16} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}.$$

Given $\mathbf{b} = [0, 2]^\top$ the solution of $A\mathbf{x} = \mathbf{b}$ is $[1, 1]^\top$, as can be easily found by Gaussian elimination. However, our computer ‘believes’ that $A^\top A$ is singular!

QR and least squares

Let A be an $m \times n$ matrix with $m \geq n$, and let $A = QR$ be a *reduced* QR factorization where Q is $m \times n$ has orthonormal columns and R is $n \times n$ upper triangular. We know that \mathbf{x} is a solution to the least squares problem iff $A\mathbf{x} - \mathbf{b}$ is orthogonal to all columns of A . Since the columns of Q span the same space as the columns of A this is equivalent to saying that $Q^\top(A\mathbf{x} - \mathbf{b}) = \mathbf{0}$. Since the columns of Q form an orthonormal system we have³ $Q^\top Q = I_n$, and so this leads to the equation $R\mathbf{x} = Q^\top \mathbf{b}$. The latter can be solved using backsubstitution.

³Note however that QQ^\top is not equal to the identity matrix! (Q is a rectangular matrix here)