

8 Lecture 8, Sep 23

Announcements

- HW2 due next Tuesday. FAQs at <http://hua-zhou.github.io/teaching/biostatm280-2016winter/biostatm280winter2016/2016/01/26/hw2-hints.html>

Last time

- LU decomposition $\mathbf{A} = \mathbf{LU}$. $2/3n^3$ flops.
- LU with partial pivoting $\mathbf{PA} = \mathbf{LU}$. $2/3n^3$ flops.
- Cholesky decomposition $\mathbf{A} = \mathbf{LL}^T$. $1/3n^3$ flops.
- Cholesky decomposition with symmetric pivoting: $\mathbf{PAP}^T = \mathbf{LL}^T$. $1/3n^3$ flops. Yields the rank of the psd matrix.
- Note the actual implementation in LAPACK may be quite different from our description.
- *No inversion mentality*. Matrix inversion takes an extra of $4/3n^3$ flops, which is unnecessary in most applications.
- Multivariate normal density evaluation.

Today

- Linear regression by Cholesky.
- QR and linear regression.
- QR by (modified) Gram-Schmidt.
- QR by Householder.
- QR by Givens.
- Sweep operator.

Linear regression by Cholesky (method of normal equations)

Assume $\mathbf{X} \in \mathbb{R}^{n \times p}$ has full column rank. (For rank deficient \mathbf{X} , use Cholesky with symmetric pivoting.)

- Cholesky on the augmented matrix $(\mathbf{X} \ \mathbf{y})^T(\mathbf{X} \ \mathbf{y})$ yields

$$\begin{pmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{y} \\ \mathbf{y}^T \mathbf{X} & \mathbf{y}^T \mathbf{y} \end{pmatrix} = \begin{pmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{l}^T & d \end{pmatrix} \begin{pmatrix} \mathbf{L}^T & \mathbf{l} \\ \mathbf{0}^T & d \end{pmatrix} = \begin{pmatrix} \mathbf{L}\mathbf{L}^T & \mathbf{L}\mathbf{l} \\ \mathbf{l}^T \mathbf{L}^T & \|\mathbf{l}\|_2^2 + d^2 \end{pmatrix}.$$

Normal equation $\mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^T \mathbf{y}$ implies the equation

$$\mathbf{L}\mathbf{L}^T \boldsymbol{\beta} = \mathbf{L}\mathbf{l} \text{ or } \mathbf{L}^T \boldsymbol{\beta} = \mathbf{l},$$

which we can solve for $\boldsymbol{\beta}$ in p^2 flops. Since $\mathbf{l} = \mathbf{L}^{-1} \mathbf{X}^T \mathbf{y}$, we have

$$\mathbf{l}^T \mathbf{l} = \mathbf{y}^T \mathbf{X} (\mathbf{L}\mathbf{L}^T)^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{y}^T \mathbf{P}_X \mathbf{y} = \|\hat{\mathbf{y}}\|_2^2$$

and

$$d^2 = \mathbf{y}^T \mathbf{y} - \mathbf{l}^T \mathbf{l} = \mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y} = \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 = \text{SSE}.$$

If standard errors are needed, we need inversion $(\mathbf{X}^T \mathbf{X})^{-1} = (\mathbf{L}\mathbf{L}^T)^{-1} = \mathbf{L}^{-T} \mathbf{L}^{-1}$. Use `chol2inv()` in R function for this purpose.

- In summary, linear regression by Cholesky, aka the method of normal equations:
 - Form the lower triangular part of $(\mathbf{X}, \mathbf{y})^T(\mathbf{X}, \mathbf{y})$.
 $n(p+1)^2 \approx np^2$ flops.
 - Cholesky decomposition of the augmented system $\begin{pmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{y} \\ \mathbf{y}^T \mathbf{X} & \mathbf{y}^T \mathbf{y} \end{pmatrix}$.
 $(p+1)^3/3 \approx p^3/3$ flops.
 - Solve $\mathbf{L}^T \boldsymbol{\beta} = \mathbf{l}$ for regression coefficients $\hat{\boldsymbol{\beta}}$.
 p^2 flops.
 - If want standard errors, estimate σ^2 by $\hat{\sigma}^2 = d^2/(n-p)$ and compute $\hat{\sigma}^2(\mathbf{X}^T \mathbf{X})^{-1} = \hat{\sigma}^2(\mathbf{L}\mathbf{L}^T)^{-1}$.
 $4p^3/3$ flops.

Total cost is $p^3/3 + np^2$ flops (without s.e.) or $5p^3/3 + np^2$ flops (with s.e.).

QR decomposition and linear regression

A second approach for linear regression uses QR decomposition. This is how the `lm()` function in R does linear regression. Assume $\mathbf{X} \in \mathbb{R}^{n \times p}$ has full column rank.

- QR decomposition: $\mathbf{X} = \mathbf{QR}$, where $\mathbf{Q} \in \mathbb{R}^{n \times n}$, $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}_n$, and $\mathbf{R} \in \mathbb{R}^{n \times p}$.
 - first p columns of \mathbf{Q} form an orthonormal basis of $\mathcal{C}(\mathbf{X})$
 - last $n - p$ columns of \mathbf{Q} form an orthonormal basis of $\mathcal{N}(\mathbf{X}^\top)$
- (Thin/Skinny QR) Then $\mathbf{X} = \mathbf{Q}_1 \mathbf{R}_1$, where $\mathbf{Q}_1 \in \mathbb{R}^{n \times p}$ has orthogonal columns, $\mathbf{Q}_1^\top \mathbf{Q}_1 = \mathbf{I}_p$, and $\mathbf{R}_1 \in \mathbb{R}^{p \times p}$ is an invertible upper triangular matrix with positive diagonal entries.
- For linear regression, we only need skinny QR.
Note $\mathbf{X}^\top \mathbf{X} = \mathbf{R}_1^\top \mathbf{R}_1$ yields the Cholesky decomposition of $\mathbf{X}^\top \mathbf{X}$.
- (Skinny) QR on the augmented matrix yields

$$\begin{pmatrix} \mathbf{X} & \mathbf{y} \end{pmatrix} = \begin{pmatrix} \mathbf{Q} & \mathbf{q} \end{pmatrix} \begin{pmatrix} \mathbf{R} & r \\ \mathbf{0}_p^\top & d \end{pmatrix} = \begin{pmatrix} \mathbf{QR} & \mathbf{Qr} + d\mathbf{q} \end{pmatrix}.$$

Normal equation $\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{y}$ implies

$$\mathbf{R}\boldsymbol{\beta} = \mathbf{R}^{-\top} \mathbf{X}^\top \mathbf{y} = \mathbf{R}^{-\top} \mathbf{R}^\top \mathbf{Q}^\top \mathbf{y} = \mathbf{Q}^\top \mathbf{y} = \mathbf{r},$$

which is easy to solve for $\boldsymbol{\beta}$. The fitted value is $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{QRR}^{-1}\mathbf{r} = \mathbf{Qr}$.
The residual is

$$\hat{\mathbf{e}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{y} - \mathbf{Qr} = d\mathbf{q}$$

and $\text{SSE} = \|\hat{\mathbf{e}}\|_2^2 = d^2$. The projection matrix is

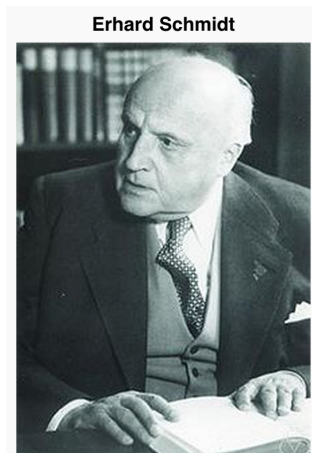
$$\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X} = \mathbf{QR}(\mathbf{R}^\top \mathbf{R})^{-1} \mathbf{R}^\top \mathbf{Q}^\top = \mathbf{QQ}^\top.$$

- Three numerical methods to compute QR: (modified) Gram-Schmidt, Householder transform, (fast) Givens transform.

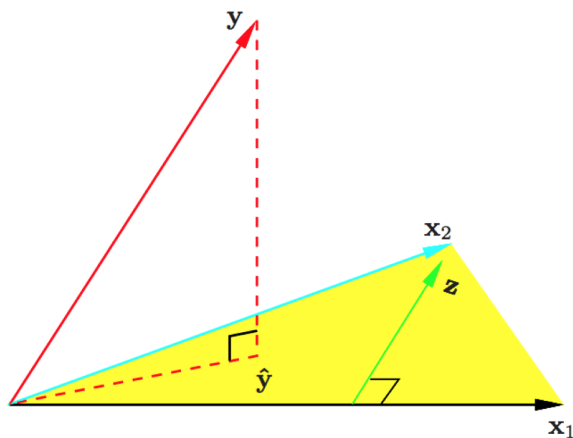
QR by (modified) Gram-Schmidt



Jørgen Pedersen Gram in an undated photo



Erhard Schmidt (courtesy MFO)



Assume $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p) \in \mathbb{R}^{n \times p}$ has full column rank.

- Gram-Schmidt (GS) algorithm produces the skinny QR: $\mathbf{X} = \mathbf{Q}\mathbf{R}$, where $\mathbf{Q} \in \mathbb{R}^{n \times p}$ has orthogonal columns and $\mathbf{R} \in \mathbb{R}^{p \times p}$ is an upper triangular matrix.
- Gram-Schmidt algorithm orthonormalizes a set of non-zero, *linearly independent* vectors $\mathbf{x}_1, \dots, \mathbf{x}_p$. Initialize $\mathbf{q}_1 = \mathbf{x}_1 / \|\mathbf{x}_1\|_2$; then for $k = 2, \dots, p$,

$$\begin{aligned} \mathbf{v}_k &= \mathbf{x}_k - \mathbf{P}_{\mathcal{C}\{\mathbf{q}_1, \dots, \mathbf{q}_{k-1}\}} \mathbf{x}_k = \mathbf{x}_k - \sum_{j=1}^{k-1} \langle \mathbf{x}_k, \mathbf{q}_j \rangle \mathbf{q}_j \\ \mathbf{q}_k &= \mathbf{v}_k / \|\mathbf{v}_k\|_2 \end{aligned}$$

- For $j = 1, \dots, p$, $\mathcal{C}\{\mathbf{x}_1, \dots, \mathbf{x}_j\} = \mathcal{C}\{\mathbf{q}_1, \dots, \mathbf{q}_j\}$, and $\mathbf{q}_j \perp \mathcal{C}\{\mathbf{x}_1, \dots, \mathbf{x}_{j-1}\}$.
- Collectively, we have $\mathbf{X} = \mathbf{Q}\mathbf{R}$ (skinny QR), where
 - $\mathbf{Q} \in \mathbb{R}^{n \times p}$ has orthonormal columns \mathbf{q}_k and thus $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}_p$.
 - Where \mathbf{R} ? $\mathbf{R} = \mathbf{Q}^T \mathbf{X} \in \mathbb{R}^{p \times p}$ has entries $r_{jk} = \langle \mathbf{q}_j, \mathbf{x}_k \rangle$, which are available from the algorithm. Note $r_{jk} = 0$ for $j > k$. Thus \mathbf{R} is upper triangular.
- In GS algorithm, \mathbf{X} is over-written by \mathbf{Q} and \mathbf{R} is stored in a separate array.
- The regular Gram-Schmidt is *unstable* (we loose orthogonality due to roundoff errors) when columns of \mathbf{X} are collinear.
- *Modified Gram-Schmidt* (MGS): after each normalization step of \mathbf{v}_k , we replace $\tilde{\mathbf{x}}_j$, $j > k$, by its residual.
- Why MGS is better than GS? Read <http://cavern.uark.edu/~arnold/4353/CGSMGS.pdf>
- Computational cost of GS and MGS is $\sum_{k=1}^p 4n(k-1) \approx 2np^2$.

QR by Householder



Photograph by [Paul Halmos](#)

Assume $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p) \in \mathbb{R}^{n \times p}$ has full column rank.

- Idea: $\mathbf{H}_p \cdots \mathbf{H}_2 \mathbf{H}_1 \mathbf{X} = \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix}$, where $\mathbf{H}_j \in \mathbb{R}^{n \times n}$ are the Householder transformation matrix. It yields the full QR where $\mathbf{Q} = \mathbf{H}_1 \cdots \mathbf{H}_p \in \mathbb{R}^{n \times n}$. Recall that GS/MGS only produces the thin QR decomposition.
- For arbitrary $\mathbf{v}, \mathbf{w} \in \mathbb{R}^n$ with $\|\mathbf{v}\|_2 = \|\mathbf{w}\|_2$, we can construct a *Householder matrix* $\mathbf{H} = \mathbf{I}_n - 2\mathbf{u}\mathbf{u}^\top$, $\mathbf{u} = -\frac{1}{\|\mathbf{v}-\mathbf{w}\|_2}(\mathbf{v} - \mathbf{w})$, that carries \mathbf{v} to \mathbf{w} :

$$\mathbf{H}\mathbf{v} = \mathbf{w}.$$

\mathbf{H} is symmetric and orthogonal. Calculation of Householder vector \mathbf{u} costs $2n$ flops.

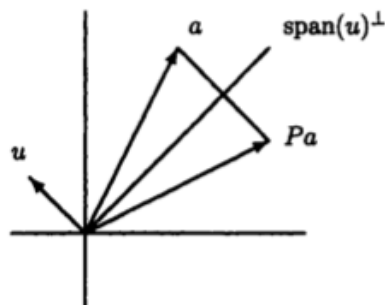


FIG. 2.3.1. Reflection of a vector \mathbf{a} in a hyperplane with normal \mathbf{u} .

- Now choose \mathbf{H}_1 to zero the first column of \mathbf{X} below diagonal

$$\mathbf{H}_1 \mathbf{x}_1 = \begin{pmatrix} \|\mathbf{x}_1\|_2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

Take \mathbf{H}_2 to zero the second column below diagonal; ...

$$\mathbf{H}_2 \mathbf{H}_1 \mathbf{A} = \begin{bmatrix} \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & 0 & \boxtimes & \times & \times \\ 0 & 0 & \boxtimes & \times & \times \\ 0 & 0 & \boxtimes & \times & \times \\ 0 & 0 & \boxtimes & \times & \times \end{bmatrix}$$

In general, choose the j -th Householder transform $\mathbf{H}_j = \mathbf{I}_n - 2\mathbf{u}_j\mathbf{u}_j^\top$, where $\mathbf{u}_j = \begin{pmatrix} \mathbf{0}_{j-1} \\ \tilde{\mathbf{u}}_j \end{pmatrix}$, $\tilde{\mathbf{u}}_j \in \mathbb{R}^{n-j+1}$, to zero the j -th column below diagonal. \mathbf{H}_j takes the form

$$\mathbf{H}_j = \begin{pmatrix} \mathbf{I}_{j-1} & \\ & \mathbf{I}_{n-j+1} - 2\tilde{\mathbf{u}}_j\tilde{\mathbf{u}}_j^\top \end{pmatrix} = \begin{pmatrix} \mathbf{I}_{j-1} & \\ & \tilde{\mathbf{H}}_j \end{pmatrix}.$$

- Applying a Householder transform $\mathbf{H} = \mathbf{I} - 2\mathbf{u}\mathbf{u}^\top$ to a matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$

$$\mathbf{H}\mathbf{X} = \mathbf{X} - 2\mathbf{u}(\mathbf{u}^\top \mathbf{X})$$

costs $4np$ flops. *We never explicitly form the Householder matrices.*

Note applying \mathbf{H}_j to \mathbf{X} only needs $4(n-j+1)(p-j+1)$ flops.

- QR by Householder: $\mathbf{H}_p \cdots \mathbf{H}_1 \mathbf{X} = \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix}$.
- The process is done in place. Upper triangular part of \mathbf{X} is overwritten by \mathbf{R}_1 and the essential Householder vectors ($\tilde{\mathbf{u}}_{j1}$ is normalized to 1) are stored in $\mathbf{X}(j:n, j)$.
- At j -th stage
 1. computing the Householder vector $\tilde{\mathbf{u}}_j$ costs $2(n-j+1)$ flops
 2. applying the Householder transform $\tilde{\mathbf{H}}_j$ to the $\mathbf{X}(j:n, j:p)$ block costs $4(n-j+1)(p-j+1)$ flops

In total we need $\sum_{j=1}^p [2(n-j+1) + 4(n-j+1)(p-j+1)] \approx 2np^2 - \frac{2}{3}p^3$ flops.

- Where is \mathbf{Q} ? $\mathbf{Q} = \mathbf{H}_1 \cdots \mathbf{H}_p$.

In some applications, it's necessary to form the orthogonal matrix \mathbf{Q} .

Accumulating \mathbf{Q} costs another $2np^2 - \frac{2}{3}p^3$ flops.

- When computing $\mathbf{Q}^\top \mathbf{v}$ or $\mathbf{Q}\mathbf{v}$ as in some applications (e.g., solve linear equation using QR), no need to form \mathbf{Q} . Simply apply Householder transforms successively.

`qr.qy()` and `qr.qty()` in R do this.

- Computational cost of Householder QR for linear regression: $2np^2 - \frac{2}{3}p^3$ (regression coefficients and $\hat{\sigma}^2$) or more (fitted values, s.e., ...).

Rank deficient \mathbf{X} : Householder QR with column pivoting

$\mathbf{X} \in \mathbb{R}^{n \times p}$ may not have full column rank.

- Idea (due to Businger and Golub 1965): at the j -th stage, swap the column in $\mathbf{X}(j : n, j : p)$ with maximum ℓ_2 norm to be the pivot column. If the maximum ℓ_2 norm is 0, it stops, ending with

$$\mathbf{X}\mathbf{\Pi} = \mathbf{Q} \begin{pmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ \mathbf{0}_{(n-r) \times r} & \mathbf{0}_{(n-r) \times (p-r)} \end{pmatrix},$$

where $\mathbf{\Pi} \in \mathbb{R}^{p \times p}$ is a permutation matrix and r is the rank of \mathbf{X} . QR with column pivoting is rank revealing.

- The overhead of re-computing the column norms can be reduced by the property

$$\mathbf{Q}\mathbf{z} = \begin{pmatrix} \alpha \\ \boldsymbol{\omega} \end{pmatrix} \Rightarrow \|\boldsymbol{\omega}\|_2^2 = \|\mathbf{z}\|_2^2 - \alpha^2$$

for any orthogonal matrix \mathbf{Q} .

- In R, the `qr()` function is a wrapper for various LINPACK (default) and LAPACK routines. It performs Householder QR with column pivoting and returns
 - `$qr`: a matrix of same size as input matrix
 - `$rank`: rank of the input matrix
 - `$pivot`: pivot vector
 - `$aux`: normalizing constants of Householder vectors

Auxiliary functions `qr.coef()`, `qr.resid()`, `qr.qy()`, `qr.qty()`, `qr.solve()`, ... are very helpful.

QR by Givens rotation



- Householder transform \mathbf{H}_j introduces batch of zeros into a vector.
- Givens transform (aka Givens rotation, Jacobi rotation, plane rotation) selectively zeros one element of a vector.
- Overall QR by Givens rotation is less efficient than the Householder method, but is better suited for matrices with structured patterns of nonzero elements.
- Givens/Jacobi rotations:

$$\mathbf{G}(i, k, \theta) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & c & s & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & -s & c & 0 \\ \vdots & \vdots & \vdots & \ddots \\ 0 & 0 & 0 & 1 \end{pmatrix},$$

where $c = \cos(\theta)$ and $s = \sin(\theta)$. $\mathbf{G}(i, k, \theta)$ is orthogonal.

- Pre-multiplication by $\mathbf{G}(i, k, \theta)^T$ rotates counterclockwise θ radians in the (i, k) coordinate plane. If $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{y} = \mathbf{G}(i, k, \theta)^T \mathbf{x}$, then

$$y_j = \begin{cases} cx_i - sx_k & j = i \\ sx_i + cx_k & j = k \\ x_j & j \neq i, k \end{cases}.$$

Apparently if we choose $\tan(\theta) = -x_k/x_i$, or equivalently,

$$c = \frac{x_i}{\sqrt{x_i^2 + x_k^2}}, \quad s = \frac{-x_k}{\sqrt{x_i^2 + x_k^2}},$$

then $y_k = 0$.

- Pre-applying Givens transform $\mathbf{G}(i, k, \theta)^T \in \mathbb{R}^{n \times n}$ to a matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ only effects two rows of \mathbf{A} :

$$\mathbf{A}([i, k], :) \leftarrow \begin{pmatrix} c & s \\ -s & c \end{pmatrix}^T \mathbf{A}([i, k], :),$$

costing $3m$ flops.

- Post-applying Givens transform $\mathbf{G}(i, k, \theta) \in \mathbb{R}^{m \times m}$ to a matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ only effects two columns of \mathbf{A} :

$$\mathbf{A}(:, [i, k]) \leftarrow \mathbf{A}(:, [i, k]) \begin{pmatrix} c & s \\ -s & c \end{pmatrix},$$

costing $3n$ flops.

- QR by Givens: $\mathbf{G}_t^T \cdots \mathbf{G}_1^T \mathbf{X} = \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0} \end{pmatrix}$

$$\begin{aligned} & \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} \xrightarrow{\langle 3,4 \rangle} \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \mathbf{0} & \times & \times \end{bmatrix} \xrightarrow{\langle 2,3 \rangle} \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \mathbf{0} & \times & \times \\ \mathbf{0} & \times & \times \end{bmatrix} \xrightarrow{\langle 1,2 \rangle} \\ & \begin{bmatrix} \times & \times & \times \\ \mathbf{0} & \times & \times \\ \mathbf{0} & \times & \times \\ \mathbf{0} & \times & \times \end{bmatrix} \xrightarrow{\langle 3,4 \rangle} \begin{bmatrix} \times & \times & \times \\ \mathbf{0} & \times & \times \\ \mathbf{0} & \times & \times \\ \mathbf{0} & \mathbf{0} & \times \end{bmatrix} \xrightarrow{\langle 2,3 \rangle} \begin{bmatrix} \times & \times & \times \\ \mathbf{0} & \times & \times \\ \mathbf{0} & \mathbf{0} & \times \\ \mathbf{0} & \mathbf{0} & \times \end{bmatrix} \xrightarrow{\langle 3,4 \rangle} \mathbf{R} \end{aligned}$$

- Zeros in \mathbf{X} can also be introduced row-by-row.
- If $\mathbf{X} \in \mathbb{R}^{n \times p}$, the total cost is $3np^2 - p^3$ flops and $O(np)$ square roots.
- Note each Givens transform can be summarized by a single number, which is stored in the zeroed entry of \mathbf{X} .
- Fast Givens transform avoids taking square roots.

Sweep operator

Assume $\mathbf{A} \succeq \mathbf{0}_{n \times n}$.

- KL 7.4-7.6; Also see “*A tutorial on the SWEEP operator*” by James H. Goodnight. <http://www.jstor.org/stable/2683825>
- *Sweep* on the k -th diagonal entry $a_{kk} \neq 0$ yields $\hat{\mathbf{A}}$ with entries

$$\begin{aligned}\hat{a}_{kk} &= -\frac{1}{a_{kk}} \\ \hat{a}_{ik} &= \frac{a_{ik}}{a_{kk}} \\ \hat{a}_{kj} &= \frac{a_{kj}}{a_{kk}} \\ \hat{a}_{ij} &= a_{ij} - \frac{a_{ik}a_{kj}}{a_{kk}}, \quad i \neq k, j \neq k.\end{aligned}$$

n^2 flops (taking into account of symmetry).

- *Inverse sweep* sends \mathbf{A} to $\check{\mathbf{A}}$ with entries

$$\begin{aligned}\check{a}_{kk} &= -\frac{1}{a_{kk}} \\ \check{a}_{ik} &= -\frac{a_{ik}}{a_{kk}} \\ \check{a}_{kj} &= -\frac{a_{kj}}{a_{kk}} \\ \check{a}_{ij} &= a_{ij} - \frac{a_{ik}a_{kj}}{a_{kk}}, \quad i \neq k, j \neq k.\end{aligned}$$

n^2 flops (taking into account of symmetry).

- $\check{\check{\mathbf{A}}} = \mathbf{A}$.
- Block form of sweep: Let the symmetric matrix \mathbf{A} be partitioned as $\mathbf{A} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix}$. If possible, sweep on the diagonal entries of \mathbf{A}_{11} yields

$$\hat{\mathbf{A}} = \begin{pmatrix} -\mathbf{A}_{11}^{-1} & \mathbf{A}_{11}^{-1}\mathbf{A}_{12} \\ \mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12} \end{pmatrix}.$$

Order dose *not* matter.

- Pd and determinant: \mathbf{A} is pd if and only if each diagonal entry can be swept in succession and is positive until it is swept. When a diagonal entry of a pd matrix \mathbf{A} is swept, it becomes negative and remains negative thereafter. Taking the product of diagonal entries just before each is swept yields the determinant of \mathbf{A} .

- Linear regression by sweep. Sweep on $\begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{y} \\ \mathbf{y}^\top \mathbf{X} & \mathbf{y}^\top \mathbf{y} \end{pmatrix}$ yields

$$\begin{pmatrix} -(\mathbf{X}^\top \mathbf{X})^{-1} & (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ \mathbf{y}^\top \mathbf{X} (\mathbf{X} \mathbf{X})^{-1} & \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \end{pmatrix} = \begin{pmatrix} -\frac{1}{\sigma^2} \text{Var}(\hat{\boldsymbol{\beta}}) & \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\beta}}^\top & \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 \end{pmatrix}.$$

In total $np^2 + p^3$ flops.

- Sweep is useful for stepwise regression, (conditional) multivariate normal density calculation, MANOVA, ...
- Warning: the `sweep()` function in R has nothing do to with the sweep operator here.
- Demo code: <http://hua-zhou.github.io/teaching/biostatm280-2016winter/sweep.html>

Summary of linear regression: Table on KL p105

Method	Flops	Remarks	Software	Stability
Sweep	$np^2 + p^3$	$(\mathbf{X}^\top \mathbf{X})^{-1}$ available	SAS	less stable
Cholesky	$np^2 + p^3/3$			less stable
QR by Householder	$2np^2 - 2/3p^3$		R	stable
QR by MGS	$2np^2$	\mathbf{Q}_1 available		more stable

Table 1: Numerical methods for linear regression. In order of stability.

Remarks:

- When $n \gg p$, sweep and Cholesky are twice faster than QR and need less space. But QR methods are more stable and produce numerically more accurate solution.

- Although sweep is slower than Cholesky, it yields standard errors and so on.
- Sweep is useful for stepwise regression, multivariate normal calculation, and numerous other statistical applications.
- MGS appears slower than Householder, but it yields \mathbf{Q}_1 .
- Sweep and Cholesky is based on the Gram matrix $\mathbf{X}^T \mathbf{X}$, which can be dynamically updated with incoming data. They can easily handle huge n , moderate p data sets that cannot fit into memory.

“There is simply no such thing as a universal ‘gold standard’ when it comes to algorithms”.