- How to compute QR factorization?

  Idea 1: convert basis to orthonormal one

  Gram – Schmidt: given $a_1, \cdots, a_n \in \mathbb{C}^m$, produce orthonormal $q_1, \cdots, q_n$ with span$\{q_1, \cdots, q_n\}$ = span$\{a_1, \cdots, a_n\}$

$$q_1 = \frac{a_1}{\|a_1\|_2}$$

$$q_2' = a_2 - (q_1^* a_2) q_1, \qquad q_2 = \frac{q_2'}{\|q_2'\|_2}$$

$$\cdots$$

$$q_j' = a_j - \sum_{i=1}^{j-1} (q_i^* a_i) q_i, \qquad q_j = \frac{q_j'}{\|q_j'\|_2} \qquad j = 3, \cdots, n$$

$\Longleftarrow$

$$a_1 = \|a_1\|_2 \, q_1$$

$$a_2 = (q_1^* a_2) q_1 + q_2 \|a_2 - (q_1^* a_2) q_1\|_2$$

$$\cdots$$

$$a_j = \sum_{i=1}^{j-1} (q_i^* a_i) q_i + q_j \|a_j - \sum_{i=1}^{j-1} (q_i^* a_i) q_i\|_2$$

$\Longleftarrow$

$$A = QR$$

$$A = [\, a_1 \cdots a_n \,], \qquad R_{ij} = \begin{cases} q_i^* a_j, & i < j \\[2mm] \|a_j - \sum_{k=1}^{j-1} (q_k^* a_k) q_k\|_2, & i = j \\[2mm] 0, & \text{otherwise} \end{cases}$$

# Implementation :

## Classical GS

Input: $A \in \mathbb{C}^{m \times n}$ of rank $n$

Output: $Q \in \mathbb{C}^{m \times n}$, $R \in \mathbb{C}^{n \times n}$

for $j = 1, \cdots, n$

    for $i = 1, \cdots, j-1$      *operation count:

        $R_{ij} = q_i^* a_j \longrightarrow (j-1)(2m-1)$

    end

    $q_j' = a_j - \sum_{i=1}^{j-1} R_{ij} q_i \qquad \longrightarrow 2m(j-1)$

    $R_{jj} = \| q_j' \|_2 , \qquad q_j = q_j' / R_{jj}$

          $\searrow 2m \qquad\qquad \searrow m$

end

Operation count:

$$\sum_{j=1}^{n} (j-1)(2m-1) + 2(j-1)m + 2m + m$$

$$\approx \sum_{j=1}^{n} 4jm$$

$$\approx 2mn^2 \text{ FLOPs}$$

- **weakness:** there is nothing to force orthonormality of $Q$ in classical Gram-Schmidt

    ex. $A = \begin{bmatrix} 1 & 1 & 1 \\ \varepsilon & & \\ & \varepsilon & \\ & & \varepsilon \end{bmatrix}$

    assume $\varepsilon$ so small such that $fl(1+\varepsilon^2) = 1$

    then $Q_{GS} = \begin{bmatrix} 1 & 0 & 0 \\ \varepsilon & -1/\sqrt{2} & -1/\sqrt{2} \\ 0 & 1/\sqrt{2} & 0 \\ 0 & 0 & 1/\sqrt{2} \end{bmatrix}$

                      $q_1 \quad q_2 \quad q_3$

        $q_2^* q_3 = \frac{1}{2}$

- **Remedy:** Instead of orthogonalizing $a_j$ at step $j$ only, can orthogonalize $a_j$ to $q_i$ as soon as $q_i$ is computed

## modified GS

Let $a_k^{(1)} = a_k$, $\quad k = 1, \cdots, n$

for $k = 1, \cdots, n$

$\quad R_{kk} = \|a_k^{(k)}\|_2$, $\quad q_k = a_k^{(k)} / R_{kk}$

$\quad$ for $j = k+1, \cdots, n$

$\quad\quad R_{kj} = q_k^* a_j^{(k)}$, $\quad a_j^{(k+1)} = a_j^{(k)} - R_{kj} q_k$

$\quad$ end

end

Classical GS is equivalent to modified GS $\longleftarrow$ exercise

<u>Thm</u> Suppose Modified GS is applied to $A \in \mathbb{R}^{m \times n}$ of rank $n$ yielding $\hat{Q} \in \mathbb{R}^{m \times n}$, $\hat{R} \in \mathbb{R}^{n \times n}$,

$\quad \exists \; c_i = c_i(m,n)$, s.t.

$\quad A + \Delta A_1 = \hat{Q}\hat{R}$, $\quad \|\Delta A_1\|_2 \leq c_1 \, \varepsilon_{mach} \|A\|_2$

$$\|\hat{Q}^T\hat{Q} - I\|_2 \leq C_2 \, \varepsilon_{mach} \, K_2(A) \Big/ \big(1 - C_2' \, \varepsilon_{mach} \, K_2(A)\big)$$

and $\exists \; Q \in \mathbb{R}^{m \times n}$ with orthonormal columns s.t.

$\hat{R}$ is the $\longrightarrow$ $A + \Delta A_2 = Q\hat{R}$, $\quad \|\Delta A_2\|_2 \leq C_3 \, \varepsilon_{mach} \|A\|_2$
exact triangular
QR factor
of a matrix $$\|Q - \hat{Q}\|_2 \leq C_4 \, \varepsilon_{mach} \, K_2(A) \Big/ \big(1 - C_4' \, \varepsilon_{mach} \, K_2(A)\big)$$
near
to A. Pf: Higham Thm 19.13
i.e. it is a good
R-factor.

the departure from orthonormality of $\hat{Q}$ is

bounded by $O(K_2(A) \, \varepsilon_{mach})$

- Back to least-squares:

  Solve least-squares via QR:

  $$\hat{x} = \hat{R}^{-1} \hat{Q}^* b$$

  use $\hat{Q}^* \hat{Q} = I_n$ suffer from large $K_2(A)$

If we translation the result to backward error analysis, we get that the computed $\hat{x}$ is the exact solution of the LS problem
$$\| b + \Delta b - (A + \Delta A) y \|_2^2$$
where $\frac{\| \Delta b \|_2}{\| b \|_2} \leq K_2(A) \, \varepsilon_{mach}$, $\frac{\| \Delta A \|_2}{\| A \|_2} \leq \varepsilon_{mach}$

NOT Backward stable

  Instead, we can apply Modified GS to $[A \ b]$

  $$[A \ b] = \begin{bmatrix} Q_1 & q_{n+1} \end{bmatrix} \begin{bmatrix} R & z \\ o & \rho \end{bmatrix}$$

  We have $\quad Ax - b = [A \ b] \begin{bmatrix} x \\ -1 \end{bmatrix}$

  $$= \begin{bmatrix} Q_1 & q_{n+1} \end{bmatrix} \begin{bmatrix} Rx - z \\ -\rho \end{bmatrix}$$

  $$= Q_1 (Rx - z) - \rho \, q_{n+1}$$

  Hence $\quad \| Ax - b \|_2^2 = \| Rx - z \|_2^2 + \rho^2$

  So $\quad x = R^{-1} z \quad$ is the Least-squares solution

  <u>Thm</u> Solving LS via Modified GS for $[A \ b]$

  has forward error as good as a backward stable algo.

---

- Can perform QR factorization faster than $2mn^2$ FLOPs

  if don't need to form Q explicitly

  * Householder:

  A

  $$\begin{pmatrix} x & x & x \\ x & x & x \\ x & x & x \\ x & x & x \end{pmatrix} \xrightarrow{Q_1} \begin{pmatrix} x & x & x \\ o & x & x \\ o & x & x \\ o & x & x \end{pmatrix} \xrightarrow[\begin{bmatrix} 1 & o \\ o & H_2 \end{bmatrix}]{Q_2} \begin{pmatrix} x & x & x \\ o & x & x \\ o & o & x \\ o & o & x \end{pmatrix} \xrightarrow[\begin{bmatrix} 1 & \\ & H_3 \end{bmatrix}]{Q_3} \begin{pmatrix} x & x & x \\ o & x & x \\ o & o & x \\ o & o & o \end{pmatrix}$$

  $H_2$

$$Q_3 Q_2 Q_1 A = \begin{bmatrix} R \\ 0 \end{bmatrix} \implies A = (Q_3 Q_2 Q_1)^* \begin{bmatrix} R \\ 0 \end{bmatrix}$$

$Q_i$ unitary

$$= Q \begin{bmatrix} R \\ 0 \end{bmatrix} \qquad \leftarrow \text{ full } QR$$

$Q \in \mathbb{C}^{m \times n}$ unitary

$\begin{bmatrix} R \\ 0 \end{bmatrix} \in \mathbb{C}^{m \times n}$ upper triangular

\* How to choose $Q_k$ ?

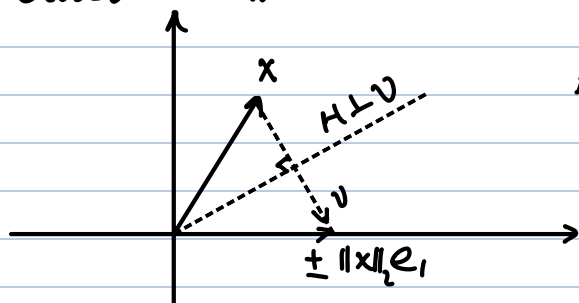$$\text{Let } Q_k = \begin{bmatrix} I & 0 \\ 0 & H_k \end{bmatrix} \begin{matrix} k-1 \\ m-(k-1) \end{matrix} \qquad , \quad H_k \text{ unitary (thus norm preserving)}$$

(over columns: $k-1 \quad m-(k-1)$)

keep the first $k-1$ coordinates unchanged

$x \in \mathbb{C}^{m-(k-1)}$ be $k$th, $\cdots$, $m$th entries of $k$th column

Goal : $H_k x = \pm \|x\|_2 e_1$

$$\begin{pmatrix} x \\ x \\ x \end{pmatrix} \xrightarrow{H_k} \begin{pmatrix} x \\ 0 \\ 0 \end{pmatrix}$$

Householder transform: reflection across the plane $H$

$$H_k x = x - 2 P_v x = x - 2\left(\frac{v^*}{\|v\|_2} x\right)\frac{v}{\|v\|_2}$$

$$= \left(I - \frac{2vv^*}{\|v\|_2^2}\right) x$$

$\underbrace{\qquad\qquad}$ Householder transform

(a) For stability, one may choose $v = \text{sign}(x_1) \|x\| e_1 + x$ to avoid catastrophic cancellation in computing $\|v\|_2^2$ when $\|x\| \approx |x_1|$

$\llcorner$ why ?

(b) Householder matrices $H_k$ is never formed in upper triangularizing $A$, storage and computations use solely the Householder vector $v$.

- Implementation:

  | Householder QR:
  |
  |  For $k = 1, \cdots, n$
  |
  |     $X \leftarrow A(k:m, k)$
  |
  |     $V_k = \text{sign}(x_{(1)}) \|x\|_2 e_1 + x$
  |
  |     $V_k \leftarrow V_k / \|V_k\|_2$
  |
  |     $A(k:m, k:n) \leftarrow A(k:m, k:n) - 2 V_k \left(V_k^* A(k:m, k:n)\right)$
  |
  | end

  Above the last equation (in blue):
  substraction $(m-k)(n-k)$ FLOPs
  outer product $(m-k)(n-k)$ FLOPs
  matrix-vector product $2(m-k)(n-k)$ FLOPs

- operation count:

$$\approx \sum_{k=1}^{n} 4(m-k)(n-k) \quad \text{FLOPs}$$

$$\approx 4\left(mn^2 - (m+n)\frac{n^2}{2} + \frac{n^3}{3}\right)$$

$$\approx 2mn^2 - \frac{2}{3}n^3$$

- Stability of Householder transform

  Since the orthogonal transform is performed by Householder vector, the orthogonality of $Q$ is enforced

  <u>Thm</u> Let $\begin{bmatrix} \hat{R} \\ 0 \end{bmatrix} \in \mathbb{R}^{m \times n}$ be the computed upper triangular QR

  factor of $A \in \mathbb{R}^{m \times n}$, Let $Q = (\hat{Q}_n \cdots \hat{Q}_1)^T$ be

  the exact orthogonal matrix obtained from Householder

  vectors computed from the algorithm

  Then $\qquad A + \Delta A = Q\hat{R} \qquad$ with $\quad \|\Delta A\|_2 \leq C_{m.n} \|A\|_2$

  Let $\hat{Q} = fl\left((\hat{Q}_n \cdots \hat{Q}_1)^T\right)$ be the computed orthogonal matrix

  Then $\qquad \hat{Q} = Q(I_m + \Delta I) \qquad$ with $\|\Delta I\|_2 \leq C_{m.n}$

  (blue annotation, arrow pointing to $\hat{Q}$:)
  $\hat{Q}$ is very close to orthogonal matrix
  regardless of $K(A)$!

  This a consequence of the backward stability of
  matrix multiplication

- Back to solve least-squares

Solve least-squares via $QR$:

$$\hat{x} = \hat{R}^{-1} \underbrace{\hat{Q}^* b}$$

use $\hat{Q}^* \hat{Q} = I_n$ NOT suffer from large $K_2(A)$

__thm__ Let $A \in \mathbb{R}^{m \times n}$ have full rank and that the least square problem $\min\limits_{x} \| Ax - b\|_2$ is solved using Householder $QR$ factorization. The computed solution $\hat{x}$ is the exact solution to

$$\min\limits_{x} \| b + \Delta b - (A + \Delta A)\hat{x} \|_2$$

where $\|\Delta A\|_2 \leq C_{min} \|A\|_2$, $\|\Delta b\|_2 \leq C_{min} \|b\|_2$

Backward stable !

Summary:

Operation count:

$O(2mn^2)$         $O(4mn^2)$

normal equation $<$ Householder $QR$ $<$ Modified GS $<$ SVD
              (without computing    (Augmented)
              $Q$ explicitly)

Rounding error:

$O(K_2^2(A) \, \varepsilon_{mach})$         $O(K_2(A) \, \varepsilon_{mach})$

normal equation $>$ Householder $QR$ $\approx$ Modified GS $\approx$ SVD
             (without computing     (Augmented)
             $Q$ explicitly)