

Michael Goforth  
CAAM 550  
HW 4  
9/22/2021

**Problem 1**  
**part a**

$$\mathbf{G}^T \mathbf{G} = \begin{bmatrix} g_1^T \\ g_2^T \\ \vdots \\ g_m^T \end{bmatrix} \begin{bmatrix} g_1 & g_2 & \dots & g_m \end{bmatrix}$$

where  $g_1, g_2, \dots, g_m$  are the columns of  $\mathbf{G}$ . Then

$$\mathbf{G}^T \mathbf{G} = \begin{bmatrix} g_1^T g_1 & g_1^T g_2 & \dots & g_1^T g_m \\ g_2^T g_1 & g_2^T g_2 & \dots & g_2^T g_m \\ \vdots & \vdots & \ddots & \vdots \\ g_m^T g_1 & g_m^T g_2 & \dots & g_m^T g_m \end{bmatrix}$$

and for a given element  $q_a, b \in \mathbf{G}^T \mathbf{G}$ , for any  $a, b \neq j, k$ ,

$$q_{a,b} = \begin{cases} 1, & a = b \\ 0, & a \neq b \end{cases}$$

Also for any element  $a \neq j, k$ ,  $q_{a,j} = 0$ ,  $q_{a,k} = 0$  and any element  $b \neq j, k$ ,  $q_{j,b} = 0$ ,  $q_{k,b} = 0$ . Finally,

$$q_{j,j} = q_{k,k} = \cos^2(\theta) + \sin^2(\theta) = 1$$

$$q_{j,k} = q_{k,j} = \cos(\theta)\sin(\theta) - \cos(\theta)\sin(\theta) = 0$$

Combining all of this together gives

$$q_{a,b} = \begin{cases} 1, & a = b \\ 0, & a \neq b \end{cases} = \mathbf{I}$$

$\mathbf{G}^T \mathbf{G} = \mathbf{I}$  so a Givens rotation is orthogonal.

**part b**

Consider a matrix  $A \in \mathbb{R}^{2 \times 2}$  where

$$\mathbf{A} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

and Givens rotation matrix

$$\mathbf{G} = \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{bmatrix}$$

If we want to use a Givens rotation to make  $A$  upper triangular, then we want

$$(GA)_{1,2} = -a \sin(\theta) + c \cos(\theta) = 0$$

$$a \sin(\theta) = c \cos(\theta)$$

$$\theta = \arctan(c/a)$$

In a similar way, we can use Givens rotations to zero out all the lower elements of a generic matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . For every element in  $a_i$  below the diagonal of  $\mathbf{A}$ , a matrix  $\mathbf{G}_j$  can be formed such that  $\mathbf{A}_j = \mathbf{G}_j \mathbf{A}$  and  $(A_j)_i = 0$ . Repeat this for every subdiagonal element and then combine the  $G$  matrices such that

$$\mathbf{G} = \prod_{i=1}^n \mathbf{G}_i$$

where each  $\mathbf{G}_i$  is a Givens rotation that zeros out one of the  $n$  subdiagonal elements of  $\mathbf{A}$ . Then

$$\mathbf{G}\mathbf{A} = (\mathbf{G}_n \dots \mathbf{G}_2 \mathbf{G}_1) \mathbf{A} = \mathbf{R}$$

which is upper triangular. Also, as shown in part a, Givens rotations are orthogonal so

$$\mathbf{G}^T \mathbf{G} \mathbf{A} = \mathbf{G}^T \mathbf{R}$$

and finally

$$\mathbf{A} = \mathbf{G}^T \mathbf{R} = \mathbf{Q} \mathbf{R}$$

One Givens rotation will be needed for each subdiagonal element. When  $n \leq m$ , this will be the  $(n-1)$ th triangular number, which can be found as  $\frac{1}{2}(n^2 - n)$ . When  $m > n$ , we will have require the same number as in the previous case plus an additional  $(m-n)$  rows of length  $n$ . Combining these yields the number of Givens rotations,  $i$ , as

$$i = \begin{cases} \frac{n}{2}(n-1), & n \leq m \\ \frac{n}{2}(n-1) + (m-n)n, & m > n \end{cases}$$

As shown above,

$$\mathbf{A} = \mathbf{G}^T \mathbf{R} = \mathbf{Q} \mathbf{R}$$

so

$$\mathbf{Q} = \mathbf{G}^T$$

### part c

See Jupyter notebook for implementation and output.

### Problem 2

$x + (y + z) = (x + y) + z$  does not hold in floating point arithmetic due to the rounding that takes place at each step. For example, if

$$\bar{x} = 5.112$$

$$\bar{y} = 5.112$$

$$\bar{z} = 5.113$$

$$\text{fl}(\text{fl}(\bar{x} + \bar{y}) + \bar{z}) = \text{fl}(\text{fl}(10.224) + 5.113) = \text{fl}(1.022 * 10^1 + 5.113) = \text{fl}(15.333) = 1.533 * 10^1$$

$$\text{fl}(\bar{x} + \text{fl}(\bar{y} + \bar{z})) = \text{fl}(5.112 + \text{fl}(10.225)) = \text{fl}(5.112 + 1.023 * 10^1) = \text{fl}(15.342) = 1.534 * 10^1$$

**Problem 3**  
**part i**

$$\begin{aligned}
\sigma &= \left[ \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2} \\
&= \left[ \frac{1}{n-1} \sum_{i=1}^n (x_i^2 - 2\bar{x}x_i + \bar{x}^2) \right]^{1/2} \\
&= \left[ \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - \sum_{i=1}^n 2\bar{x}x_i + \sum_{i=1}^n \bar{x}^2 \right) \right]^{1/2} \\
&= \left[ \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - 2\bar{x} \frac{n}{n} \sum_{i=1}^n x_i + n\bar{x}^2 \right) \right]^{1/2} \\
&= \left[ \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - 2\bar{x}n\bar{x} + n\bar{x}^2 \right) \right]^{1/2} \\
&= \left[ \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) \right]^{1/2}
\end{aligned}$$

**part ii** See Jupyter notebook for code and results.

**Problem 4**  
**part i**

$$\begin{aligned}
A1 &= \text{fl}(4 * \text{fl}(\text{fl}(\pi) * \text{fl}(\text{fl}(r) * \text{fl}(r)))) \\
&= \text{fl}(4 * \text{fl}(\text{fl}(3.142) * \text{fl}(\text{fl}(6370) * \text{fl}(6370)))) \\
&= \text{fl}(4 * \text{fl}(3.142 * \text{fl}(6.370 * 10^3 * 6.370 * 10^3))) \\
&= \text{fl}(4 * \text{fl}(3.142 * \text{fl}(40,576,900))) \\
&= \text{fl}(4 * \text{fl}(3.142 * 4.058 * 10^7)) \\
&= \text{fl}(4 * \text{fl}(127,502,360)) \\
&= \text{fl}(4 * 1.275 * 10^8) \\
&= \text{fl}(510,000,000) \\
&= 5.100 * 10^8
\end{aligned}$$

**part ii**

$$\begin{aligned}
A2 &= \text{fl}(4 * \text{fl}(\text{fl}(\pi) * \text{fl}(\text{fl}(r) * \text{fl}(r)))) \\
&= \text{fl}(4 * \text{fl}(\text{fl}(3.142) * \text{fl}(\text{fl}(6371) * \text{fl}(6371)))) \\
&= \text{fl}(4 * \text{fl}(3.142 * \text{fl}(6.371 * 10^3 * 6.371 * 10^3))) \\
&= \text{fl}(4 * \text{fl}(3.142 * \text{fl}(40,589,641))) \\
&= \text{fl}(4 * \text{fl}(3.142 * 4.059 * 10^7)) \\
&= \text{fl}(4 * \text{fl}(127,533,780)) \\
&= \text{fl}(4 * 1.275 * 10^8) \\
&= \text{fl}(510,000,000) \\
&= 5.100 * 10^8
\end{aligned}$$

$$A2 - A1 = \text{fl}(5.100 * 10^8 - 5.100 * 10^8) = \text{fl}(0) = 0$$

**part iii**

$$\begin{aligned}
\frac{d}{dr}A &= \text{fl}(8\text{fl}(\pi r)) \\
&= \text{fl}(8\text{fl}(3.142 * 6.370 * 10^3)) \\
&= \text{fl}(8\text{fl}(20,014.54)) \\
&= \text{fl}(8 * 2.001 * 10^4) \\
&= \text{fl}(160,080) \\
&= 1.601 * 10^5
\end{aligned}$$

**part iv**

Part iii is more accurate.

**part v**

Part iii is more accurate because it is computing the difference directly. Part i and ii are the same in floating point arithmetic because the magnitude of the areas is much greater than the difference between them, so when the result of part i is subtracted from the result of part ii, catastrophic cancellation occurs. If floating point arithmetic with a longer mantissa is used, the difference between parts i and ii will be more accurate.