

[09-Sep-2021] Homework 1

Please work the following problems from our text:

- Problems 2.4, 2.5, 2.9, 2.11
- Problems 3.1, 3.3, 3.11, 3.12

2.4 Suppose $f : \mathbb{R} \rightarrow \mathbb{R}$ is infinitely differentiable, and we wish to write algorithms for finding x^* minimizing $f(x)$. Our algorithm outputs x_{est} , an approximation of x^* . Assuming that in our context this problem is equivalent to finding roots of $f'(x)$, write expressions for:

- Forward error of the approximation
- Backward error of the approximation
- Conditioning of this minimization problem near x^*

This is the same as solving $g(x) = 0$ where $g(x) = f'(x)$ and the analysis is the same as Example 2.10. Setting $x^* = x + \varepsilon$ and replacing f by g yields

Example 2.10 (Root-finding). Suppose that we are given a smooth function $g : \mathbb{R} \rightarrow \mathbb{R}$ and want to find roots x with $g(x) = 0$. By Taylor's theorem, $g(x + \varepsilon) \approx g(x) + \varepsilon g'(x)$ when $|\varepsilon|$ is small. Thus, an approximation of the condition number for finding the root x is given by

$$\frac{\text{forward error}}{\text{backward error}} = \frac{|(x + \varepsilon) - x|}{|g(x + \varepsilon) - g(x)|} \stackrel{(a)}{\approx} \frac{|\varepsilon|}{|\varepsilon g'(x)|} \stackrel{(b)}{\approx} \frac{1}{|g'(x)|} = \frac{1}{|f''(x)|} \stackrel{(c)}{\approx}$$

This approximation generalizes the one in Example 2.9. If we do not know x , we cannot evaluate $g'(x)$, but if we can examine the form of g and bound $|g'|$ near x , we have an idea of the worst-case situation.

- 2.5 Suppose we are given a list of floating-point values x_1, x_2, \dots, x_n . The following quantity, known as their “log-sum-exp,” appears in many machine learning algorithms:

$$\ell(x_1, \dots, x_n) \equiv \ln \left[\sum_{k=1}^n e^{x_k} \right].$$

- (a) The value $p_k \equiv e^{x_k}$ often represents a probability $p_k \in (0, 1]$. In this case, what is the range of possible x_k 's?
- (b) Suppose many of the x_k 's are very negative ($x_k \ll 0$). Explain why evaluating the log-sum-exp formula as written above may cause numerical error in this case.
- (c) Show that for any $a \in \mathbb{R}$,

$$\ell(x_1, \dots, x_n) = a + \ln \left[\sum_{k=1}^n e^{x_k - a} \right].$$

To avoid the issues you explained in 2.5b, suggest a value of a that may improve the stability of computing $\ell(x_1, \dots, x_n)$.

(a) $\{x : e^x \in (0, 1]\} \approx \{\log p : p \in (0, 1]\} = (-\infty, 0)$

(b) If $x_k \ll 0$ for many values of k , then some of the terms in the sum are very small compared to the others. When a small number is added to a large number, the small number loses precision.

In the present case evaluating $\ln(\alpha+\beta)$ where α is large compared to β is best done as

$$\ln(\alpha+\beta) = \ln \left(\alpha \left(1 + \frac{\beta}{\alpha} \right) \right) = \ln \alpha + \ln \left(1 + \frac{\beta}{\alpha} \right) = \ln \alpha + \text{log1p} \left(\frac{\beta}{\alpha} \right)$$

where $\text{log1p}(\varepsilon) = \ln(1+\varepsilon)$ is evaluated using the built-in function obtained from the Taylor series

$$\ln(1+\varepsilon) = \varepsilon - \frac{1}{2}\varepsilon^2 + \frac{1}{3}\varepsilon^3 - \frac{1}{4}\varepsilon^4 + \dots \quad \text{for } \varepsilon \text{ very small.}$$

(c) $\ln \left[\sum_{k=1}^n e^{x_k} \right]$ may be evaluated more accurately by dividing the sum into two sums: a sum of very small terms plus a sum of larger terms. In particular, let $\mathcal{S} = \{k : x_k \ll 0\}$ and $\mathcal{X} = \{1, \dots, n\} \setminus \mathcal{S}$. Thus

$$\sum_{k=1}^n e^{x_k} = \alpha + \beta \quad \text{where } \alpha = \sum_{k \in \mathcal{X}} e^{x_k} \text{ and } \beta = \sum_{k \in \mathcal{S}} e^{x_k}$$

It follows, setting $a = \ln \alpha$ that $\alpha = e^a$

$$\ln \sum_{k=1}^n e^{x_k} = \ln \alpha + \ln \left(1 + \frac{\beta}{\alpha} \right) = a + \ln \left(1 + \beta e^{-a} \right)$$

$$= a + \ln \left(\left(\frac{\alpha}{\alpha + \beta} e^{-a} \right)^{e^{\alpha a}} \right) = a + \ln \left(\sum_{k=1}^n e^{x_k - a} \right)$$

Note, in this case it is still important to evaluate $\ln \left(\sum_{k=1}^n e^{x_k - a} \right)$

using the log1p function as $\text{log1p}\left(\frac{\beta}{\alpha}\right)$ as suggested earlier.

2.9 In this problem, we continue to explore the conditioning of root-finding. Suppose $f(x)$ and $p(x)$ are smooth functions of $x \in \mathbb{R}$.

- (a) Thanks to inaccuracies in how we evaluate or express $f(x)$, we might accidentally compute roots of a perturbation $f(x) + \varepsilon p(x)$. Take x^* to be a root of f , so $f(x^*) = 0$. If $f'(x^*) \neq 0$, for small ε we can write a function $x(\varepsilon)$ such that $f(x(\varepsilon)) + \varepsilon p(x(\varepsilon)) = 0$, with $x(0) = x^*$. Assuming such a function exists and is differentiable, show:

$$\frac{dx}{d\varepsilon} \Big|_{\varepsilon=0} = -\frac{p(x^*)}{f'(x^*)}.$$

- (b) Assume $f(x)$ is given by Wilkinson's polynomial [131]:

$$f(x) \equiv (x-1) \cdot (x-2) \cdot (x-3) \cdots \cdot (x-20).$$

We could have expanded $f(x)$ in the monomial basis as $f(x) = a_0 + a_1 x + a_2 x^2 + \cdots + a_{20} x^{20}$, for appropriate choices of a_0, \dots, a_{20} . If we express the coefficient a_{19} inaccurately, we could use the model from Exercise 2.9a with $p(x) \equiv x^{19}$ to predict how much root-finding will suffer. For these choices of $f(x)$ and $p(x)$, show:

$$\frac{dx}{d\varepsilon} \Big|_{\varepsilon=0, x^*=j} = -\prod_{k \neq j} \frac{j}{j-k}.$$

- (c) Compare $\frac{dx}{d\varepsilon}$ from the previous part for $x^* = 1$ and $x^* = 20$. Which root is more stable to this perturbation?

(a) Use implicit differentiation

$$\frac{d}{d\varepsilon} (f(x(\varepsilon)) + \varepsilon p(x(\varepsilon))) = 0,$$

Thus $f'(x(\varepsilon)) x'(\varepsilon) + p(x(\varepsilon)) + \varepsilon p'(x(\varepsilon)) x'(\varepsilon) = 0$

Set $\varepsilon=0$ to obtain:

$$f'(x^*) x'(0) + p(x^*) = 0$$

Therefore

$$\frac{dx}{d\varepsilon} \Big|_{\varepsilon=0} = x'(0) = -\frac{p(x^*)}{f'(x^*)}$$

- (b) Set $f(x) \equiv (x-1) \cdot (x-2) \cdot (x-3) \cdots \cdot (x-20)$.
 $p(x) \equiv x^{19}$

Then for $x^* = j$ we have

$$\frac{dx}{d\varepsilon} \Big|_{\varepsilon=0} = -\frac{p(j)}{f'(j)}$$

Differentiating,

$$f'(x) \approx \frac{d}{dx} \sum_{k=1}^n \pi_k(x-k) \approx \sum_{l=1}^k \pi_l(x-k)$$

so

$$f'(j) = \pi_{k \neq j}(j-k)$$

It follows that

$$\left. \frac{dx_k}{de} \right|_{e=0} = \frac{-\pi(j)}{f'(j)} = \frac{-j^{19}}{\pi_{k \neq j}(j-k)} = -\pi \left(\frac{j}{j-k} \right)$$

(c) When $j=1$ we obtain

$$\left. \frac{dx_k}{de} \right|_{e=0} = -\pi \left(\frac{1}{1-k} \right) \approx 3.5477$$

```
julia> ks=[1:19;]
19-element Vector{Int64}:
 1
 2
 3
 4
 5
 6
 7
 8
 9
10
11
12
13
14
15
16
17
18
19
20

julia> sum(20 ./ (20 .. -ks))
70.95479314287364
```

```
julia> ks=[2:20;]
19-element Vector{Int64}:
 2
 3
 4
 5
 6
 7
 8
 9
10
11
12
13
14
15
16
17
18
19
20

julia> sum(1 ./ (1 .. -ks))
-3.5477396571436826
```

when $j=20$ we obtain

$$\left. \frac{dx_k}{de} \right|_{e=0} = -\pi \left(\frac{20}{20-k} \right) \approx 70.955$$

It follows that the root at $j=1$ is more stable.

- 2.11 One technique for tracking uncertainty in a calculation is the use of *interval arithmetic*. In this system, an uncertain value for a variable x is represented as the interval $[x] \equiv [\underline{x}, \bar{x}]$ representing the range of possible values for x , from \underline{x} to \bar{x} . Assuming infinite-precision arithmetic, give update rules for the following in terms of \underline{x} , \bar{x} , \underline{y} , and \bar{y} :

- $[x] + [y]$
- $[x] - [y]$
- $[x] \times [y]$
- $[x] \div [y]$
- $[x]^{1/2}$

Additionally, propose a conservative modification for finite-precision arithmetic.

- $[\underline{x}] + [\underline{y}] = [\underline{\underline{x}} + \underline{y}, \bar{x} + \bar{y}]$

$\begin{matrix} \text{lower two bounds} \\ \text{is a lower bound} \end{matrix}$ $\begin{matrix} \text{upper bounds} \\ \text{is a upper bound.} \end{matrix}$

$\begin{matrix} \text{mult by } \underline{y} \text{ changes} \\ \text{order of inequalities} \end{matrix}$

- $[\underline{x}] - [\bar{y}] = [\underline{x}, \bar{x}] - [\underline{y}, \bar{y}] = [\underline{\underline{x}}, \bar{x}] + [-\bar{y}, \underline{y}] = [\underline{\underline{x}} - \bar{y}, \bar{x} - \underline{y}]$

- $[\underline{x}] \times [\underline{y}]$ There are lots of cases depending on whether \underline{x} , \bar{x} , \underline{y} and \bar{y} are positive or negative.

One solution is to define

$$\Delta = \{ \underline{x}\underline{y}, \underline{x}\bar{y}, \bar{x}\underline{y}, \bar{x}\bar{y} \}$$

and then set

$$[\underline{x}] \times [\underline{y}] = [\min \Delta, \max \Delta]$$

- $[\underline{x}] \div [\underline{y}]$ For this we can use multiplication after taking the reciprocal of $[\underline{y}]$.

Note the special case in which

$$0 \in [\underline{y}, \bar{y}]$$

should be treated as divide by zero.

Thus $\frac{[\underline{x}]}{[\underline{y}]} = \begin{cases} \text{undet} & \text{in } 0 \in [\underline{y}, \bar{y}] \\ \left[\min \left\{ \frac{1}{\underline{y}}, \frac{1}{\bar{y}} \right\}, \max \left\{ \frac{1}{\underline{y}}, \frac{1}{\bar{y}} \right\} \right] & \text{otherwise} \end{cases}$

Then define $[\underline{x}] \div [\underline{y}] = [\underline{x}] \times \frac{1}{[\underline{y}]}$.

- $\sqrt{[x]}$ when $x \leq 0$ this might be the square root of a negative number, which is imaginary. While it is possible to define rules for the propagation of errors in complex arithmetic, that's outside the scope of the present problem.

Therefore

$$\sqrt{[x]} = \begin{cases} \text{undef} & \text{if } x \leq 0 \\ [\sqrt{\underline{x}}, \sqrt{\bar{x}}] & \text{otherwise} \end{cases}$$

In the case of finite precision arithmetic, it's useful to consider two additional rounding functions.

$\text{rdown}(x)$ and $\text{rup}(x)$ for $x \in \mathbb{R}$
which respectively round x up or down. Specifically

let $F = \{x_* : x_* \text{ is a representable floating point number}\}$

$$\text{rdown}(x) = \max \{x_* \in F : x_* \leq x\}$$

$$\text{rup}(x) = \min \{x_* \in F : x_* \geq x\}.$$

After noting that these rounding modes are available on all modern processors, we define

- $[x] + [y] = [\text{rdown}(x+y), \text{rup}(x+y)]$
- $[x] - [y] = [\text{rdown}(x-y), \text{rup}(x-y)]$
- $[x] \times [y] = [\text{rdown}(\min A), \text{rup}(\max A)]$

Note that

$$rdown(\min A) = \min \{rdown(\underline{x}\underline{y}), rdown(\underline{x}\bar{y}), rdown(\bar{x}\underline{y}), rdown(\bar{x}\bar{y})\}$$

so there is no difficulty computing this quantity on the computer

- $\lceil x \rceil \div \lceil y \rceil$ define

$$\bullet \frac{\lceil x \rceil}{\lceil y \rceil} = \begin{cases} \text{undif} & \text{if } 0 \in [\underline{y}, \bar{y}] \\ \left[\min \left\{ rdown \frac{1}{\underline{y}}, rdown \frac{1}{\bar{y}} \right\}, \max \left\{ rup \frac{1}{\underline{y}}, rup \frac{1}{\bar{y}} \right\} \right] & \text{otherwise} \end{cases}$$

and then compute $\lceil x \rceil \times \frac{1}{\lceil y \rceil}$ using the previous rule.

- $\widetilde{\lceil x \rceil}$

It may not be the case that round up and round down are available for the square root

In fact it's possible the subroutine implementing $\sqrt(x)$ does not even follow round to nearest over the entire range of available values for x .

In this case suppose y_* represents the approximation computed by the $\sqrt(x)$ function and the relative error of this approximation is bounded as

$$\left| \frac{\sqrt(x) - \sqrt{x}}{\sqrt{x}} \right| \leq \varepsilon \quad \text{for all } x \in F.$$

$$\text{Then } -\varepsilon \sqrt{x} \leq \sqrt(x) - \sqrt{x} \leq \varepsilon \sqrt{x}$$

$$\text{implies } (1-\varepsilon)\sqrt{x} \leq \sqrt(x) \leq (1+\varepsilon)\sqrt{x}$$

or that $\frac{\sqrt{x}}{1+\varepsilon} \leq \sqrt{x} \leq \frac{\sqrt{x}}{1-\varepsilon}$

consequently

$$\text{rdown}\left(\frac{\sqrt{x}}{1+\varepsilon}\right) \leq \sqrt{x} \leq \text{rup}\left(\frac{\sqrt{x}}{1-\varepsilon}\right)$$

note that everything here and here
can be calculated on the computer.

Thus

$$\widetilde{\sqrt{x}} = \begin{cases} \text{undef} & \text{if } x < 1 \\ \left[\text{rdown}\left(\frac{\sqrt{x}}{1+\varepsilon}\right), \text{rup}\left(\frac{\sqrt{x}}{1-\varepsilon}\right) \right] & \text{otherwise} \end{cases}$$

Note that other bounds may also be possible for $\widetilde{\sqrt{x}}$

3.1 Can all matrices $A \in \mathbb{R}^{n \times n}$ be factored $A = LU$? Why or why not?

No, some matrices require row swaps in order to be factored. For example

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

Can not be factored into LU because $a_{11}=0$.

DH 3.3 Factor the following matrix A as a product $A = LU$:

$$\begin{pmatrix} 1 & 2 & 7 \\ 3 & 5 & -1 \\ 6 & 1 & 4 \end{pmatrix}.$$

$$\begin{bmatrix} 1 & 2 & 7 \\ 3 & 5 & -1 \\ 6 & 1 & 4 \end{bmatrix}$$

$$r_2 \leftarrow r_2 - 3r_1$$
$$r_3 \leftarrow r_3 - 6r_1$$

$$\begin{bmatrix} 1 & 2 & 7 \\ 0 & -1 & -22 \\ 0 & -11 & -38 \end{bmatrix}$$

$$r_3 \leftarrow r_3 - 11r_2$$

$$\begin{array}{r} 22 \\ 11 \\ \hline 22 \\ 22 \\ \hline 248 \\ -38 \\ \hline 204 \end{array}$$

$$U = \begin{bmatrix} 1 & 2 & 7 \\ 0 & -1 & -22 \\ 0 & 0 & 204 \end{bmatrix}$$

$$L = \begin{bmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ 6 & 11 & 1 \end{bmatrix}$$

```
julia> U=[1 2 7; 0 -1 -22; 0 0 204]
3x3 Matrix{Int64}:
 1  2   7
 0 -1  -22
 0  0  204

julia> L=[1 0 0; 3 1 0; 6 11 1]
3x3 Matrix{Int64}:
 1  0   0
 3  1   0
 6  11  1

julia> L*U
3x3 Matrix{Int64}:
 1  2   7
 3  5  -1
 6  1   4

julia>
```

Checks...

- 3.11 Show how the LU factorization of $A \in \mathbb{R}^{n \times n}$ can be used to compute the determinant of A .

The cofactor expression for determinant is

$$\det A = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det A_{ij}$$

where A_{ij} is the matrix formed by crossing out the i th row and j th column of the matrix A .

Suppose U is an upper triangular matrix. Then

$$U = \begin{bmatrix} u_{11} & u_{12} & \cdots & u_{1n} \\ 0 & u_{22} & \cdots & u_{2n} \\ 0 & 0 & \ddots & \vdots \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & u_{nn} \end{bmatrix}$$

In this case $u_{ij} = 0$ for $i > j$.

It follows taking $i=1$ that

$$\det U = \sum_{j=1}^n (-1)^{1+j} u_{1j} \det U_{1j} = u_{11} \det U_{11}$$

Note that U_{11} is

$$U_{11} = \begin{bmatrix} u_{11} & u_{12} & \cdots & u_{1n} \\ 0 & u_{22} & \cdots & u_{2n} \\ 0 & 0 & \ddots & \vdots \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & u_{nn} \end{bmatrix}$$

also an upper triangular matrix...

Therefore $\det U_{11}$ can be computed as $u_{22} \det$

$$\begin{bmatrix} u_{33} & \cdots & u_{3n} \\ 0 & \ddots & \vdots \\ 0 & \cdots & 0 & u_{nn} \end{bmatrix}$$

and so on so $\det U = u_{11} u_{22} \cdots u_{nn}$ is the product along the diagonal.

A similar argument implies

$$\det L = 1$$

where

$$L = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ l_{21} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ l_{n1} & \cdots & \cdots & 1 \end{bmatrix}$$

since there are only 1's on the diagonal.

Thus,

$$\det(A) = \det(LU) = \det(L)\det(U) = u_1u_2 \cdots u_n$$

Therefore to compute the determinant of A, first find the LU factorization and then multiply the diagonal elements of U together to get the answer.

- 3.12 For numerical stability and generality, we incorporated pivoting into our methods for Gaussian elimination. We can modify our construction of the LU factorization somewhat to incorporate pivoting as well.

- Argue that following the steps of Gaussian elimination on a matrix $A \in \mathbb{R}^{n \times n}$ with partial pivoting can be used to write $U = L_{n-1}P_{n-1} \cdots L_2P_2L_1P_1A$, where the P_i 's are permutation matrices, the L_i 's are lower triangular, and U is upper triangular.
- Show that P_i is a permutation matrix that swaps rows i and j for some $j \geq i$. Also, argue that L_i is the product of matrices of the form $I_{n \times n} + c\vec{e}_k\vec{e}_i^\top$ where $k > i$.
- Suppose $j, k > i$. Show $P_{jk}(I_{n \times n} + c\vec{e}_k\vec{e}_i^\top) = (I_{n \times n} + c\vec{e}_j\vec{e}_i^\top)P_{jk}$, where P_{jk} is a permutation matrix swapping rows j and k .
- Combine the previous two parts to show that
$$L_{n-1}P_{n-1} \cdots L_2P_2L_1P_1 = L_{n-1}L'_{n-2}L'_{n-3} \cdots L'_1P_{n-1} \cdots P_2P_1,$$
where L'_1, \dots, L'_{n-2} are lower triangular.
- Conclude that $A = PLU$, where P is a permutation matrix, L is lower triangular, and U is upper triangular.
- Extend the method from §3.5.2 for solving $A\vec{x} = \vec{b}$ when we have factored $A = PLU$, without affecting the time complexity compared to factorizations $A = LU$.

(4) Elimination with partial pivoting consists of a sequence of elimination steps of the form $r_i \leftarrow r_i - \alpha r_j$ where $i > j$ followed by row swaps $r_i \leftrightarrow r_j$.

(b) Since $A = IA_j$, the elimination step $r_i \leftarrow r_i - \alpha r_j$ corresponds to multiplication by a matrix formed by performing the elimination step on the identity matrix. In particular, the matrix for

$$[r_i \leftarrow r_i - \alpha r_j] = I - \alpha e_i e_j^\top$$

is lower triangular since $i > j$ and the matrix

$$[r_i \leftrightarrow r_j] = I - e_i e_i^\top + e_j e_j^\top - e_j e_i^\top + e_i e_j^\top$$

is a permutation matrix that swaps row i with row j .

The reduction to row echelon form results in the upper triangular matrix U obtained as

$$[r_0 \leftarrow r_n - \alpha_{n,n} r_{n-1}] \dots [r_2 \leftarrow r_k] [r_n \leftarrow r_n - \alpha_{n,k} r_k] \dots [r_i \leftarrow r_i - \alpha_{i,k} r_k] A = U$$

Since the inverse of $r_i \leftarrow r_i - \alpha r_j$ is $r_i \leftarrow r_i + \alpha r_j$ and also triangular, each of these matrices can be inverted in order and placed on the other side to obtain

$$A = [r_i \leftarrow r_j] [r_2 \leftarrow r_2 + \alpha_{2,i} r_i] \dots [r_n \leftarrow r_n + \alpha_{n,i} r_i] [r_2 \leftarrow r_n] \dots [r_n \leftarrow r_n + \alpha_{n,n-1} r_{n-1}] U.$$

P_1 L_1 elimination steps between row swap P_2 \dots L_{n-1}

Now group the elimination steps between each row swap together

Since the product of lower triangular matrices is again lower triangular than each of the L_i 's are lower triangular. We obtain

$$A = P_1 L_1 P_2 \dots L_{n-1} U$$

(c) If $j, k > i$ then $j \neq i$ and $k \neq i$. Now, since

performing a row swap followed by an elimination step is the same as performing the elimination step on the row that was swapped first and then swapping it; then the sequence

$$r_k \leftarrow r_k + c r_i \quad \text{followed by} \quad r_j \leftarrow r_k$$

is the same as

$$r_j \leftarrow r_k \quad \text{followed by} \quad r_j \leftarrow r_j + c r_i$$

Written in matrix notation, this means

$$P_{j,k} (I + C e_k e_i^\top) = (I + C e_j e_i^\top) P_{j,k}$$

Note that the order of the operations represented by the matrix multiplication read from right to left as matrix multiplication really means composition of linear functions.

(d) One can switch the order of the L and P matrices using the result of (c) to group all the L matrices together.

Thus $A = P_1 L_1 P_2 \dots L_{n-1} U = P_1 P_2 \dots P_{n-1} L_1' L_2' \dots L_{n-1}' U$
against lower triangular

(e) Since multiplication of lower triangular matrices are again lower triangular than $L \equiv L_1' L_2' \dots L_{n-1}'$ is lower triangular.

Since $P = P_1 P_2 \dots P_{n-1}$ is also a permutation, then

$$A = PLU.$$

(f) To solve $Ax = b$ use the fact that $P^T = P^{-1}$ and simplify $PLUx = b$ to obtain $LUX = P^T b$.

At this point one can solve

$$\begin{cases} Ly = P^T b \\ Ux = y \end{cases}$$

as before using forward and back substitution.