



## DDA 3005 — Numerical Methods

### Solutions 2

**Problem 1 (Evaluating Functions and Condition Number):** (approx. 25 points)

In this exercise, we consider the mathematical problem of evaluating functions  $f : \mathbb{R} \rightarrow \mathbb{R}$ .

- a) Show that the problem  $f(x) = \log(1+x)$  is well-conditioned for  $x \geq -\frac{1}{2}$ .

**Hint:** The estimate  $\log(1+x) \leq x$  (for all  $x > -1$ ) might be useful.

- b) Compute the relative condition number  $\text{cond}_f(x)$  of  $f(x) := \frac{\sin(x)}{x}$ . Is the evaluation of  $f$  at  $x = 0$  a well-conditioned problem or not? Is the evaluation of  $f$  generally well-conditioned (for all  $x$ )? Provide detailed explanations!

**Solution :**

- a) The condition number is given by

$$\text{cond}(x) = \left| \frac{x f'(x)}{f(x)} \right| = \left| \frac{x}{(1+x) \log(1+x)} \right| =: g(x)$$

The function  $g$  can be shown to be differentiable (for  $x > -1$ ) and we have

$$g(x) = \begin{cases} \frac{x}{(1+x) \log(1+x)} & \text{if } x \neq 0, \\ 0 & \text{if } x = 0, \end{cases} \quad \text{and} \quad g'(x) = \begin{cases} \frac{\log(1+x)-x}{(1+x)^2 \log(1+x)^2} & \text{if } x \neq 0, \\ -\frac{1}{2} & \text{if } x = 0. \end{cases}$$

(The limits for  $x = 0$  can be verified by, e.g., applying L'Hôpital). Using the hint, this implies  $g'(x) \leq 0$  and hence,  $g$  is monotonically decreasing. This yields  $\text{cond}(x) \leq g(-\frac{1}{2}) \approx 1.44 \leq 2$ . Thus, evaluating the function  $f$  is well-conditioned.

- b) The condition number for  $f(x) = \frac{\sin(x)}{x}$  is given by

$$\text{cond}(x) = \left| \frac{x f'(x)}{f(x)} \right| = \left| \frac{x^2 \cdot \frac{\cos(x)x - \sin(x)}{x^2}}{\sin(x)} \right| = \left| \frac{\cos(x)x}{\sin(x)} - 1 \right| =: g(x)$$

Using L'Hôpital, we have  $\lim_{x \rightarrow 0} \frac{\cos(x)x}{\sin(x)} = \lim_{x \rightarrow 0} \frac{-\sin(x)x + \cos(x)}{\cos(x)} = 1$  and hence, it follows  $g(0) = 0$ . Consequently, evaluating  $f$  at  $x = 0$  is “perfectly conditioned”. The evaluation of  $f$  is generally not well-conditioned for all  $x$ .

For  $x \rightarrow k\pi$ ,  $k \in \mathbb{Z} \setminus \{0\}$ , we obtain  $g(x) \rightarrow \infty$ . Thus, evaluating  $f$  is not well-conditioned for such  $x$ .

**Problem 2 (Clustering and Flops):**

(approx. 25 points)

Let  $\mathbf{x}^1, \dots, \mathbf{x}^m \in \mathbb{R}^n$  be a given point cloud of  $m$  different vectors in  $\mathbb{R}^n$ . In this exercise, we want to study an algorithm that allows to cluster the points  $\mathbf{x}^1, \dots, \mathbf{x}^m$  into different groups or clusters. Here, clustering refers to assigning a given vector  $\mathbf{x}^i$  to a specific group or cluster that  $\mathbf{x}^i$  belongs to. We label the groups  $1, \dots, k$  and specify a clustering or assignment of the  $m$  vectors to groups using an  $m$ -dimensional vector  $\mathbf{c}$ , where  $c_i$  is the group number that  $\mathbf{x}^i$  is assigned to. With each of the groups  $1, \dots, k$ , we associate a group representative  $\mathbf{z}^1, \dots, \mathbf{z}^k \in \mathbb{R}^n$ . Naturally, the (not necessarily known) group representative should be close to the vectors in its associated group, i.e., the terms  $\|\mathbf{x}^i - \mathbf{z}^{c_i}\|$  should be small. Our goal is to find a choice of group assignments  $c_1, \dots, c_m$  and a choice of representatives  $\mathbf{z}^1, \dots, \mathbf{z}^k$  that minimize the clustering objective

$$f(\mathbf{c}, \mathbf{z}^1, \dots, \mathbf{z}^k) := \frac{1}{m} \left( \|\mathbf{x}^1 - \mathbf{z}^{c_1}\|^2 + \dots + \|\mathbf{x}^m - \mathbf{z}^{c_m}\|^2 \right).$$

We want to use a heuristic algorithm to estimate  $\mathbf{c}$  and the group representatives  $\mathbf{z}^1, \dots, \mathbf{z}^k$ . The full procedure is shown below.

Given:  $k \in \mathbb{N}$ ; a list of  $m$  vectors  $\mathbf{x}^1, \dots, \mathbf{x}^m$ ; and an initial list of  $k$  group representative vectors  $\mathbf{z}^1, \dots, \mathbf{z}^k$ . Repeat until STOP:

1. Partition the vectors into  $k$  groups. For each vector  $i = 1, \dots, m$  assign  $\mathbf{x}^i$  to the group associated with the nearest representative  $\mathbf{z}^j$ , i.e., find  $j$  such that  $\|\mathbf{x}^i - \mathbf{z}^j\|$  is minimal and set  $c_i = j$ .
2. Update representatives. For each group  $j = 1, \dots, k$ , set  $\mathbf{z}^j$  to be the mean of the vectors in the group  $j$ .

- a) Estimate the total number of flops this algorithm requires per iteration. Express your final result in terms of asymptotic big- $\mathcal{O}$  notation.

**Hint:** You can assume that the Euclidean norm is used when calculating distances. The square root operation  $\sqrt{\cdot}$  can be realized using 1 flop per application. Furthermore, finding the minimum element  $\min\{a_1, \dots, a_n\}$  of an array of  $n$  numbers costs  $\mathcal{O}(n)$  flops (the total costs are bounded by  $2n - 2$  flops).

- b) The codes `clustering.m` and `clustering.py` implement the clustering algorithm in **MATLAB** and **Python**. Download the code from BB and write a test program that allows to verify your result from part b) experimentally. Specifically, generate three random data clouds  $\mathbf{A} = [\mathbf{a}^1, \dots, \mathbf{a}^p]$ ,  $\mathbf{B} = [\mathbf{b}^1, \dots, \mathbf{b}^p]$ , and  $\mathbf{C} = [\mathbf{c}^1, \dots, \mathbf{c}^p]$  where

$$\mathbf{a}^i = \begin{bmatrix} r_1 \\ r_2 \end{bmatrix}, \quad \mathbf{b}^i = 4 \begin{bmatrix} r_3 \\ r_4 \end{bmatrix} - \begin{bmatrix} 4 \\ 6 \end{bmatrix}, \quad \mathbf{c}^i = 2 \begin{bmatrix} r_5 \\ r_6 \end{bmatrix} - \begin{bmatrix} 5 \\ 2 \end{bmatrix}, \quad r_j \sim \mathcal{N}(0, 1)$$

for all  $j = 1, \dots, 6$  and  $i = 1, \dots, p$  and set  $\mathbf{X} = [\mathbf{x}^1, \dots, \mathbf{x}^m] = [\mathbf{A}, \mathbf{B}, \mathbf{C}]$  as test set. Run the code for different data sizes  $p = 100 \cdot 1:10$  (or  $p = 100 \cdot 1:5$ ) and report the mean cpu-time per iteration for each run (e.g., using a plot). Do the observed results match your expectation? Visualize the obtained clustering results for the largest choice of  $p$ .

**Hint:** You can set the number of groups to 3 and initialize  $\mathbf{z}^1, \mathbf{z}^2, \mathbf{z}^3$  randomly. The number of iterations per run can be controlled via the `options`; you can use 100 or 1000.

**Solution :**

- a) For each  $i = 1, \dots, m$ , we first compute the distances:

$$\|\mathbf{x}^i - \mathbf{z}^1\|, \quad \dots, \quad \|\mathbf{x}^i - \mathbf{z}^k\|.$$

Calculating  $\|\mathbf{x}^i - \mathbf{z}^j\|$  costs  $n + (2n - 1) + 1 = 3n$  flops ( $n$  flops for computing  $\zeta^j = \mathbf{x}^i - \mathbf{z}^j$ ;  $2n - 1$  flops for the inner product  $(\zeta^j)^\top \zeta^j$ ; 1 flop for taking the square root. Consequently, this first step requires  $3nk$  flops in total. As mentioned in the hint, the operation  $c_i = \arg \min_{j=1, \dots, k} \|\mathbf{x}^i - \mathbf{z}^j\|$  costs additional  $2k - 2$  flops. Hence, in order to determine the full group assignment  $\mathbf{c}$ , we require  $(3nk + 2k - 2) \cdot m$  flops.

For all  $j = 1, \dots, k$ , we now set  $\mathcal{I}_j := \{i : c_i = j\}$ . The new group representative  $\mathbf{z}^j$  is then calculated via  $\mathbf{z}^j = \frac{1}{|\mathcal{I}_j|} \sum_{i \in \mathcal{I}_j} \mathbf{x}^{c_i}$ . This operation costs  $(|\mathcal{I}_j| - 1) \cdot n + 1$  flops (for the summations and final multiplication). Hence, in total, the second step of the algorithm requires  $\sum_{j=1}^k (|\mathcal{I}_j| - 1) \cdot n + 1 = (m - k)n + k$  flops.

In total this yields  $3kmn + 2km + mn - kn + k - 2m = \mathcal{O}(kmn)$  flops.

- b) The MATLAB code is shown below (to save space, the Python code is provided in the corresponding zip-file on BB).

```

1 % demo-clustering
2 rng(20240929)
3 tests = 10;
4 t      = zeros(tests,1);
5
6 for fac = 1:tests
7     A      = norminv(rand(fac*100,2), 0, 1)';
8     B      = 4*norminv(rand(fac*100,2), 0, 1)' - [4;6];
9     C      = 2*norminv(rand(fac*100,2), 0, 1)' - [5;2];
10
11     X      = [A,B,C];
12     Z0     = [A(:,1),B(:,1),C(:,1)];
13
14     opts.iter = 1000;
15
16     tic
17     [c,Z]    = clustering(X,Z0,opts);
18     t(fac) = toc/opts.iter;
19 end
20
21 figure
22 plot(1:tests,t,'.-','MarkerSize',10,'LineWidth',1.2);
23 set(gcf,'Renderer','painters');
24 saveas(gcf,strcat('exercise_24_01'),'eps');
25
26 figure
27 hold on
28 plot(X(1,c==1),X(2,c==1),'.','Color',[25,50,108]/255,'MarkerSize',12);
29 plot(X(1,c==2),X(2,c==2),'.','Color',[223,149,86]/255,'MarkerSize',12);
30 plot(X(1,c==3),X(2,c==3),'.','Color',[53,130,134]/255,'MarkerSize',12);
31
32 plot(Z(1,:),Z(2,:), 's','MarkerSize',10,'MarkerFaceColor',[1,0,0],'MarkerEdgeColor',
    ,[1,0,0]);
33 hold off

```

```

34 set(gcf, 'Renderer', 'painters');
35 saveas(gcf, strcat('exercise_24_02'), 'eps');

```

The runtime and clustering result are shown below in Figure 1. Based on the complexity result derived in part a), we expect the runtime to grow linearly with the number of datapoints  $m$  (or with  $p$ ). This is also reflected in Figure 1 (a).

**Remark:** The algorithm described and analyzed in Problem 2 is known as *k-means*. *k-means* is of fundamental importance in machine learning and clustering.

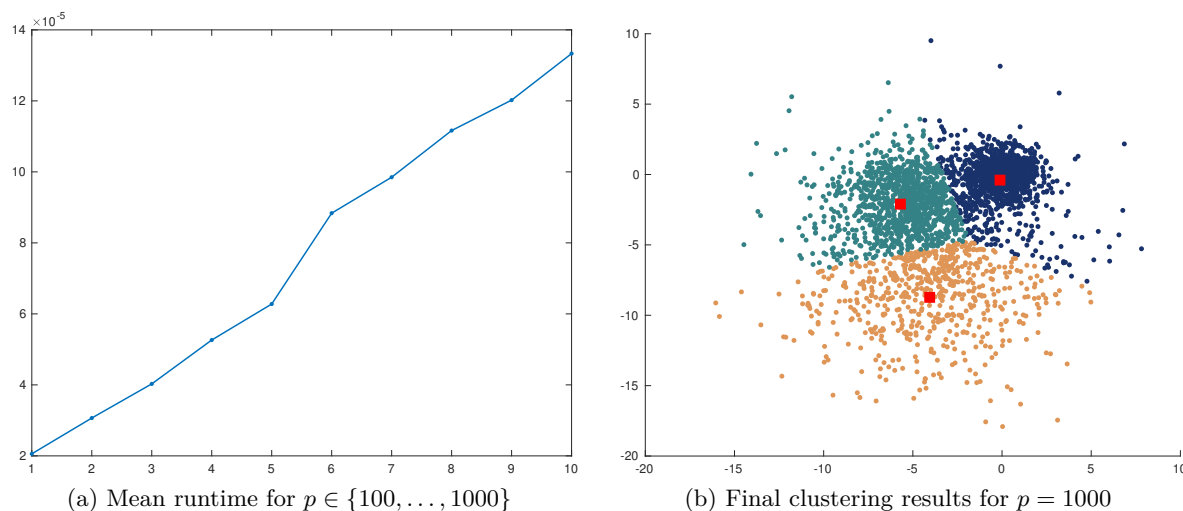


Figure 1: Results and visualization for Problem 2.

### Problem 3 (Big- $\mathcal{O}$ Calculus):

(approx. 10 points)

In this exercise, we want to verify computational rules involving the asymptotic big- $\mathcal{O}$  notation.

- Show that  $(1 + \mathcal{O}(x))(1 + \mathcal{O}(x)) = 1 + \mathcal{O}(x)$  for  $x \rightarrow 0$ .

The precise meaning of this statement is that if  $f$  is a function satisfying  $f(x) = (1 + \mathcal{O}(x))(1 + \mathcal{O}(x))$  as  $x \rightarrow 0$ , then  $f$  also satisfies  $f(x) = 1 + \mathcal{O}(x)$  as  $x \rightarrow 0$ .

**Solution :** Per definition of the asymptotic big- $\mathcal{O}$  notation, there exist functions  $f(x) = \mathcal{O}(x)$  and  $g(x) = \mathcal{O}(x)$  and constants  $C_f$ ,  $C_g$ , and  $\varepsilon$  such that  $|f(x)| \leq C_f \cdot |x|$  and  $|g(x)| \leq C_g \cdot |x|$  for all  $x$  with  $|x| \leq \varepsilon$ . We then have

$$(1 + f(x)) \cdot (1 + g(x)) = 1 + \underbrace{f(x) + g(x) + f(x) \cdot g(x)}_{=:h(x)}.$$

To establish the desired result, we need to show  $h(x) = \mathcal{O}(x)$  for  $x \rightarrow 0$ . For all  $x$  with  $|x| \leq \varepsilon$ , we have

$$\begin{aligned} h(x) &= f(x) + g(x) + f(x) \cdot g(x) \leq |f(x)| + |g(x)| + |f(x)| \cdot |g(x)| \\ &\leq C_f|x| + C_g|x| + C_f C_g |x|^2 \leq (C_f + C_g + C_f C_g \varepsilon) \cdot |x|. \end{aligned}$$

Hence, it follows  $h(x) = \mathcal{O}(x)$  as  $x \rightarrow 0$  which verifies  $(1 + \mathcal{O}(x))(1 + \mathcal{O}(x)) = 1 + \mathcal{O}(x)$  for  $x \rightarrow 0$ .

**Problem 4 (Error Analysis):**

(approx. 30 points)

In this exercise, we consider the problem of computing the function value:

$$F(\mathbf{z}) = \begin{bmatrix} x^2 - y^2 \\ 2xy \end{bmatrix}, \quad \mathbf{z} = \begin{bmatrix} x \\ y \end{bmatrix} \in \mathbb{R}^2, \quad \mathbf{z} \neq \mathbf{0}. \quad (1)$$

The following algorithm is used to evaluate  $F$ :

$$\hat{F}(\mathbf{z}) = \begin{bmatrix} (\text{fl}(x) \odot \text{fl}(x)) \ominus (\text{fl}(y) \odot \text{fl}(y)) \\ 2 \odot (\text{fl}(x) \odot \text{fl}(y)) \end{bmatrix}.$$

Here, we use a base 2 floating-point system with machine precision  $\varepsilon_{\text{mach}}$  which satisfies:

- $|\text{fl}(x) - x| \leq \varepsilon_{\text{mach}}|x|$  for all  $x \in \mathbb{R}$ .
- The IEEE-Standard 754 holds, i.e., for any machine numbers  $x, y$ , we have  $x \circledast y = \text{fl}(x \ast y)$ , where  $\ast$  can represent the operations  $\{+, -, \cdot, /\}$ .

The goal of this problem is to investigate accuracy of the algorithm  $\hat{F}$ .

- Show that the algorithm  $\hat{F}$  is accurate for problem (1). You can use the following steps:
  - First express  $\hat{F}(\mathbf{z})$  in terms of the true inputs  $x, y \in \mathbb{R}$  and  $\varepsilon_{\text{mach}}$  (using the asymptotic Landau big- $\mathcal{O}$  notation).
 

**Hint:** The following fact can be useful. Let  $\varepsilon_i, i = 1, \dots, m$  be a collection of numbers satisfying  $|\varepsilon_i| \leq \varepsilon_{\text{mach}}$ . Then, it holds that  $\prod_{i=1}^m (1 + \varepsilon_i) = 1 + \mathcal{O}(\varepsilon_{\text{mach}})$ .
  - Estimate the relative forward error  $\|\hat{F}(\mathbf{z}) - F(\mathbf{z})\|/\|F(\mathbf{z})\|$ .
- Is  $\hat{F}$  also a stable algorithm for evaluating  $F$ ? Explain your answer briefly!

**Solution :**

- We first start with part i). According to the assumption, there are  $|\varepsilon_i| \leq \varepsilon_{\text{mach}}, i = 1, 2, \dots, 7$ , such that  $\text{fl}(x) = (1 + \varepsilon_1)x, \text{fl}(y) = (1 + \varepsilon_3)y$ ,

$$\begin{aligned} & (\text{fl}(x) \odot \text{fl}(x)) \ominus (\text{fl}(y) \odot \text{fl}(y)) \\ &= \{[x(1 + \varepsilon_1)]^2(1 + \varepsilon_2) - [y(1 + \varepsilon_3)]^2(1 + \varepsilon_4)\}(1 + \varepsilon_5) \\ &= x^2(1 + \varepsilon_1)^2(1 + \varepsilon_2)(1 + \varepsilon_5) - y^2(1 + \varepsilon_3)^2(1 + \varepsilon_4)(1 + \varepsilon_5) \end{aligned}$$

and

$$2 \odot (\text{fl}(x) \odot \text{fl}(y)) = 2xy(1 + \varepsilon_1)(1 + \varepsilon_3)(1 + \varepsilon_6)(1 + \varepsilon_7).$$

Thus, using the hint, we have

$$\hat{F}(\mathbf{z}) = \begin{pmatrix} x^2(1 + \mathcal{O}(\varepsilon_{\text{mach}})) - y^2(1 + \mathcal{O}(\varepsilon_{\text{mach}})) \\ 2xy(1 + \mathcal{O}(\varepsilon_{\text{mach}})) \end{pmatrix}$$

We now show that this algorithm is accurate (part ii):

$$\begin{aligned} \frac{\|\hat{F}(\mathbf{z}) - F(\mathbf{z})\|^2}{\|F(\mathbf{z})\|^2} &= \frac{(x^2 \cdot \mathcal{O}(\varepsilon_{\text{mach}}) - y^2 \cdot \mathcal{O}(\varepsilon_{\text{mach}}))^2 + 4x^2y^2 \cdot \mathcal{O}(\varepsilon_{\text{mach}})^2}{(x^2 - y^2)^2 + 4x^2y^2} \\ &\leq \frac{x^4 \cdot \mathcal{O}(\varepsilon_{\text{mach}}^2) - 2x^2y^2 \cdot \mathcal{O}(\varepsilon_{\text{mach}})^2 + y^4 \cdot \mathcal{O}(\varepsilon_{\text{mach}}^2)}{(x^2 + y^2)^2} + \mathcal{O}(\varepsilon_{\text{mach}}^2) \\ &\leq \mathcal{O}(\varepsilon_{\text{mach}}^2), \end{aligned}$$

where we used  $(x^2 - y^2)^2 + 4x^2y^2 \geq 4x^2y^2$  and  $x^4 \cdot \mathcal{O}(\varepsilon_{\text{mach}}^2) - 2x^2y^2 \cdot \mathcal{O}(\varepsilon_{\text{mach}})^2 + y^4 \cdot \mathcal{O}(\varepsilon_{\text{mach}}^2) \leq (x^2 + y^2)^2 \cdot \mathcal{O}(\varepsilon_{\text{mach}}^2)$ . Taking the square root, this implies that  $\hat{F}$  is accurate.

- b) Notice that accuracy implies stability (we can simply choose  $\hat{z} = z$ ). Hence,  $\hat{F}$  is also a stable algorithm!

**Problem 5 (Nonexistent LU Factorization):**

(approx. 10 points)

Prove that the matrix

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 1 \end{bmatrix}$$

is nonsingular and that it does not possess a (naïve) LU factorization (without pivoting).

**Solution :** Using the cofactor expansion of the determinant, we can calculate

$$\det(\mathbf{A}) = 0 \cdot \det \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix} - 1 \cdot \det \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} + 0 \cdot \det \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix} = -1 \neq 0.$$

Hence,  $\mathbf{A}$  is nonsingular. To show that  $\mathbf{A}$  has no (naïve) LU factorization, we can use the ansatz

$$\mathbf{A} = \mathbf{LU} = \begin{bmatrix} 1 & & \\ \ell_{21} & 1 & \\ \ell_{31} & \ell_{32} & 1 \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ & u_{22} & u_{23} \\ & & u_{33} \end{bmatrix} = \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ \ell_{21}u_{11} & \ell_{21}u_{12} + u_{22} & \ell_{21}u_{13} + u_{23} \\ * & * & * \end{bmatrix}.$$

This readily implies  $u_{11} = 0$ ,  $u_{12} = 1$ , and  $u_{13} = 0$  which leads to the contradiction  $0 = \ell_{21}u_{11} = 1$ . Consequently,  $\mathbf{A}$  cannot have a classical LU factorization (without pivoting).