

# Numerical Analysis: Midterm

(30 marks, only the 3 best questions count)

Urbain Vaes

October 24, 2022

**Question 1** (Floating point arithmetic, 10 marks). True or false? (+1/0/-1)

1. Let  $(\bullet)_2$  denote binary representation. It holds that  $(0.1011)_2 + (0.0101)_2 = 1$ .
2. Let  $(\bullet)_3$  denote base 3 representation. It holds that  $(1000)_3 \times (0.002)_3 = 2$ .
3. A natural number with binary representation  $(b_4b_3b_2b_1b_0)_2$  is even if and only if  $b_0 = 0$ .
4. In Julia, `Float64(.4) == Float32(.4)` evaluates to `true`.
5. Machine addition  $\hat{+}$  is a commutative operation. More precisely, given any two double-precision floating point numbers  $x \in \mathbf{F}_{64}$  and  $y \in \mathbf{F}_{64}$ , it holds that  $x \hat{+} y = y \hat{+} x$ .
6. Let  $\mathbf{F}_{32}$  and  $\mathbf{F}_{64}$  denote respectively the sets of single and double precision floating point numbers. It holds that  $\mathbf{F}_{32} \subset \mathbf{F}_{64}$ .
7. The machine epsilon of a floating point format is the smallest strictly positive number that can be represented exactly in the format.
8. Let  $\mathbf{F}_{64}$  denote the set of double precision floating point numbers. For any  $x \in \mathbf{R}$  such that  $x \in \mathbf{F}_{64}$ , it holds that  $x + 1 \in \mathbf{F}_{64}$ .
9. Let  $a_i \in \{0, 1\}$  for  $i \in \{1, 2, 3\}$ . If  $(a_1a_2a_3)_2$  is a multiple of 3, then  $(a_1a_2a_3)_4$  is a multiple of 6. Here  $(\bullet)_4$  denotes base 4 representation.
10. Let  $f: \mathbf{R} \rightarrow \mathbf{R}$  denote the function that maps  $x \in \mathbf{R}$  to the number of double precision floating point numbers contained in the interval  $[x - 1, x + 1]$ . Then  $f$  is a decreasing function of  $x$ .
11. Let  $n \in \mathbf{N}$ . The number of bits in the binary representation of  $n$  is less than or equal to 4 times the number of digits in the decimal representation of  $n$ .
12. It holds that  $(0.\overline{2200})_3 = (0.9)_{10}$ .
13. Let  $p \in \mathbf{N}$ . The set  $\{(b_0.b_1b_2 \dots b_{p-1})_2 : b_i \in \{0, 1\}\}$  contains  $2^p$  distinct real numbers.

**Question 2** (Interpolation and approximation, 10 marks). Throughout this exercise, we assume that  $x_0 < \dots < x_n$  are distinct values and that  $u: \mathbf{R} \rightarrow \mathbf{R}$  is a smooth function. The notation  $\mathbf{P}(n)$  denotes the set of polynomials of degree less than or equal to  $n$ .

1. (4 marks) Are the following statements true or false? (+1/0/-1)

- There exists a unique polynomial  $p \in \mathbf{P}(n)$  such that

$$\forall i \in \{0, \dots, n\}, \quad p(x_i) = u(x_i). \quad (1)$$

- Assume that  $p \in \mathbf{P}(n)$  is such that (1) is satisfied. Then there is a constant  $K \in \mathbf{R}$  independent of  $x$  such that

$$\forall x \in \mathbf{R}, \quad u(x) - p(x) = K(x - x_0) \dots (x - x_n).$$

- Assume that  $p \in \mathbf{P}(n)$  is such that (1) is satisfied. Then  $p$  is of degree exactly  $n$ .
- If  $x_0, \dots, x_n$  are the roots of the Chebyshev polynomial of degree  $n$ , then

$$\sup_{x \in \mathbf{R}} |(x - x_0) \dots (x - x_n)| \leq \frac{\pi}{2^n}.$$

- The function  $S: \mathbf{N} \rightarrow \mathbf{R}$  given by

$$S(n) = \sum_{i=1}^n (i + i^2 + i^3 + i^4)$$

is a polynomial of degree 5. (More precisely, there exists a polynomial of degree 5, say  $q$ , such that  $S(n) = q(n)$  for all  $n \in \mathbf{N}$ .)

2. For  $i \in \{0, \dots, n\}$ , let  $u_i = u(x_i)$ , and let  $m \leq n$  be a given natural number. We wish to fit the data  $(x_0, u_0), \dots, (x_n, u_n)$  with a function  $\hat{u}: \mathbf{R} \rightarrow \mathbf{R}$  of the form

$$\hat{u}(x) = \alpha_0 + \alpha_1 x + \dots + \alpha_m x^m.$$

Specifically, we wish to find coefficients  $\boldsymbol{\alpha} = (\alpha_0, \dots, \alpha_m)^T$  such that the error

$$J(\boldsymbol{\alpha}) := \frac{1}{2} \sum_{i=0}^n |u_i - \hat{u}(x_i)|^2$$

is minimized. Throughout this exercise, we use the notations

$$A \begin{pmatrix} 1 & x_0 & \dots & x_0^m \\ \vdots & \vdots & & \vdots \\ 1 & x_n & \dots & x_n^m \end{pmatrix}, \quad \mathbf{b} := \begin{pmatrix} u_0 \\ \vdots \\ u_n \end{pmatrix}$$

- **(3 marks)** Show that  $J(\boldsymbol{\alpha})$  may be rewritten as

$$J(\boldsymbol{\alpha}) = \frac{1}{2}(\mathbf{A}\boldsymbol{\alpha} - \mathbf{b})^T(\mathbf{A}\boldsymbol{\alpha} - \mathbf{b}).$$

- **(2 marks)** Prove that if  $\boldsymbol{\alpha}_* \in \mathbf{R}^{m+1}$  is a minimizer of  $J$ , then

$$\mathbf{A}^T \mathbf{A} \boldsymbol{\alpha}_* = \mathbf{A}^T \mathbf{b}. \tag{2}$$

- **(1 mark)** Find a solution to (2) in terms of  $u_0, \dots, u_n$  and  $n$  when  $m = 0$ . Explain.

**Question 3** (Numerical integration, 10 marks). The Gauss–Legendre quadrature formula with  $n$  nodes is an approximate integration formula of the form

$$I(u) := \int_{-1}^1 u(x) dx \approx \sum_{i=1}^n w_i u(x_i) =: \widehat{I}_n(u), \quad (3)$$

which is exact when  $u$  is a polynomial of degree less than or equal to  $2n - 1$ . (Note that the nodes are here numbered starting from 1.)

1. (5 marks) Find the nodes and weights of the Gauss–Legendre rule with  $n = 3$  nodes.
2. (2 marks) Let  $\{L_0, L_1, \dots\}$  denote orthogonal polynomials for the inner product

$$\langle f, g \rangle := \int_{-1}^1 f(x)g(x) dx$$

which, in addition, satisfy the following two conditions:

- For all  $i \in \mathbf{N}$ , the polynomial  $L_i$  is of degree  $i$ .
- The leading coefficient of  $L_i$ , which multiplies  $x^i$ , is equal to 1.

Calculate  $L_0, L_1, L_2$  and  $L_3$ . What is the connection between  $L_3$  and the rule found in the first item?

3. Assume that  $x_1, \dots, x_n$  and  $w_1, \dots, w_n$  are such that (3) is satisfied for all  $u \in \mathbf{P}(2n-1)$ .
  - (2 marks) Show that the weights are given by

$$\forall i \in \{1, \dots, n\}, \quad w_i = \int_{-1}^1 \ell_i(x) dx,$$

where  $\ell_i$  is the Lagrange polynomial

$$\ell_i(x) = \prod_{j \neq i} \frac{x - x_j}{x_i - x_j}.$$

- (1 marks) Show that the weights are all positive:  $w_i > 0$  for all  $i$ .

4. (Bonus +2) Prove the following error estimate: if  $u$  is a smooth function, then

$$|I(u) - \widehat{I}_n(u)| \leq \frac{C_{2n}}{(2n)!} \int_{-1}^1 (L_n(x))^2 dx, \quad C_{2n} := \sup_{\xi \in [-1,1]} |u^{(2n)}(\xi)|.$$

**Hint:** You may find it useful to proceed as follows:

- First show that

$$I(u) - \widehat{I}_n(u) = \int_{-1}^1 u(x) - p(x) dx, \quad (4)$$

for *any* polynomial  $p \in \mathbf{P}(2n - 1)$  such that

$$\forall i \in \{1, \dots, n\}, \quad p(x_i) = u(x_i).$$

- Notice that equation (4) is true in particular when  $p$  is the Hermite interpolation of  $u$  at the nodes  $x_1, \dots, x_n$ . Finally, conclude by using the formula for the interpolation error proved in class: if  $p$  is the Hermite interpolant of  $u$  at the nodes  $x_1, \dots, x_n$ , then

$$\forall x \in \mathbf{R}, \quad u(x) - p(x) = \frac{u^{(2n)}(\xi(x))}{(2n)!} (x - x_1)^2 \dots (x - x_n)^2.$$

**Question 4** (Vector and matrix norms, 10 marks). The 1-norm and the  $\infty$ -norm of a vector  $\mathbf{x} \in \mathbf{R}^n$  are defined as follows:

$$\|\mathbf{x}\|_1 = |x_1| + \cdots + |x_n| \quad \text{and} \quad \|\mathbf{x}\|_\infty = \max\{|x_1|, \dots, |x_n|\}.$$

These norms both induce a matrix norm through the formula

$$\|\mathbf{A}\|_p := \sup\{\|\mathbf{A}\mathbf{x}\|_p : \|\mathbf{x}\|_p = 1\}.$$

Prove, for  $\mathbf{A} \in \mathbf{R}^{n \times n}$ , that

- (10 marks)**  $\|\mathbf{A}\|_1$  is given by the maximum absolute column sum:

$$\|\mathbf{A}\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|. \quad (5)$$

- (Bonus +2)**  $\|\mathbf{A}\|_\infty$  is given by the maximum absolute row sum:

$$\|\mathbf{A}\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|.$$

**Hint:** In order to prove (5), you may find it useful to proceed as follows:

- Introduce  $j_*$  as the index of the column with maximum absolute sum:

$$j_* = \arg \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|.$$

- Prove the direction  $\geq$  in (5) by finding a vector  $\mathbf{x}$  with  $\|\mathbf{x}\|_1 = 1$  such that

$$\|\mathbf{A}\mathbf{x}\|_1 = \sum_{i=1}^n |a_{ij_*}|.$$

- Prove the direction  $\leq$  in (5) by showing that, for any  $\mathbf{x} \in \mathbf{R}^n$  with  $\|\mathbf{x}\|_1 = 1$ ,

$$\|\mathbf{A}\mathbf{x}\|_1 \leq \sum_{i=1}^n |a_{ij_*}|.$$