

Numerical Analysis: Midterm

(30 marks, only the 3 best questions count)

Urbain Vaes

December 2, 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.

Solution. The correct answers are the following:

1. True
2. True
3. True
4. False, because the binary representation of 0.4 is infinite.
5. True, because

$$x \hat{+} y = \text{fl}(x + y) = \text{fl}(y + x) = y \hat{+} x,$$

where fl is the rounding operator.

6. True
7. False. The smallest number that can be represented in a format is $2^{E_{\min} - (p-1)}$, and the machine epsilon is $2^{-(p-1)}$.
8. False, otherwise there would be infinitely many numbers in the set \mathbf{F}_{64} .
9. False. For example, $(110)_2 = 6$ and $(110)_4 = 20$.
10. False since $\lim_{x \rightarrow -\infty} f(x) = 0$ and $f(0) > 0$.
11. True. Indeed, let d denote the number of digits in the decimal representation of n . Then $n \leq 10^d - 1$. With $4d$ bits, all the numbers up to $2^{4d} - 1$ can be represented, and since $2^{4d} - 1 = 16^d - 1 \geq 10^d - 1$, the statement is true.
12. True because

$$(0.\overline{2200})_3 = (0.2200)_3 \left(1 + 3^{-4} + (3^{-4})^2 + (3^{-4})^3 + \cdots \right) = \left(\frac{2}{3} + \frac{2}{9} \right) \frac{1}{1 - 3^{-4}} = \frac{8}{9} \frac{81}{80} = \frac{9}{10}.$$

13. True because there are 2^p choices for the bits, and distinct sets of bits correspond to 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

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

- 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 $n \in \mathbf{N}$.)

Solution. The correct answers are the following:

- True. Indeed assume that p and q both satisfy (1). Then $p - q \in \mathbf{P}(n)$ and

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

Therefore $p - q$ has at least $n + 1$ roots which, given that $p - q$ is of degree at most n , is possible only if $p - q = 0$.

- False, because if it were true, then it would hold that

$$u(x) = p(x) + K(x - x_0) \dots (x - x_n),$$

implying that u is a polynomial of degree $n + 1$. Therefore, the equation cannot be true for a general smooth function u .

- False. The statement is not true in general since, if (for example) u is the function everywhere equal to zero, then the only $p \in \mathbf{P}(n)$ that satisfies (1) is $p = 0$, which is not a polynomial of degree exactly n .
- False, because the supremum on the left-hand side is equal to ∞ as

$$\lim_{x \rightarrow \infty} |(x - x_0) \dots (x - x_n)| = \infty.$$

- True.

△

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

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

- (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.

Solution.

- Notice that

$$\mathbf{A}\boldsymbol{\alpha} = \begin{pmatrix} \alpha_0 + \alpha_1 x_0 + \cdots + \alpha_m x_0^m \\ \vdots \\ \alpha_0 + \alpha_1 x_n + \cdots + \alpha_m x_n^m \end{pmatrix} = \begin{pmatrix} \widehat{u}(x_0) \\ \vdots \\ \widehat{u}(x_n) \end{pmatrix}.$$

Therefore

$$\frac{1}{2} \sum_{i=1}^n |\widehat{u}(x_i) - u_i|^2 = \frac{1}{2} \sum_{i=1}^n |(\mathbf{A}\boldsymbol{\alpha} - \mathbf{b})_i|^2 = \frac{1}{2} (\mathbf{A}\boldsymbol{\alpha} - \mathbf{b})^T (\mathbf{A}\boldsymbol{\alpha} - \mathbf{b})$$

- A necessary condition is that $\nabla J(\boldsymbol{\alpha}_*) = 0$. We calculate that

$$\frac{\partial}{\partial x_i} (\mathbf{b}^T \mathbf{x}) = \frac{\partial}{\partial x_i} \left(\sum_{j=1}^n b_j x_j \right) = \sum_{j=1}^n b_j \delta_{ij} = b_i.$$

Similarly, for any matrix $\mathbf{M} \in \mathbf{R}^{n \times n}$, it holds that

$$\frac{\partial}{\partial x_i} (\mathbf{x}^T \mathbf{M} \mathbf{x}) = \frac{\partial}{\partial x_i} \left(\sum_{j=1}^n \sum_{k=1}^n m_{jk} x_j x_k \right) = \sum_{j=1}^n \sum_{k=1}^n m_{jk} \frac{\partial}{\partial x_i} (x_j x_k).$$

Applying the formula for the derivative of a product, we obtain

$$\begin{aligned} \frac{\partial}{\partial x_i} (\mathbf{x}^T \mathbf{M} \mathbf{x}) &= \sum_{j=1}^n \sum_{k=1}^n m_{jk} \delta_{ij} x_k + m_{jk} x_j \delta_{ik} \\ &= \sum_{k=1}^n m_{ik} x_k + \sum_{j=1}^n m_{ji} x_j = (\mathbf{M} \mathbf{x} + \mathbf{M}^T \mathbf{x})_i. \end{aligned}$$

Employing these formulae, we calculate that (representing the gradient with a column vector)

$$\nabla_{\boldsymbol{\alpha}} (\mathbf{b}^T \boldsymbol{\alpha}) = \mathbf{b}, \quad \nabla_{\boldsymbol{\alpha}} (\boldsymbol{\alpha}^T \mathbf{A}^T \mathbf{A} \boldsymbol{\alpha}) = 2\mathbf{A}^T \mathbf{A} \boldsymbol{\alpha}.$$

It is then simple to conclude.

- In this case $\mathbf{A}^T \mathbf{A} = n + 1$ and α_* is a scalar. The solution is given by

$$\alpha_* = \frac{u_0 + \cdots + u_n}{n + 1},$$

which is the average of the values u_0, \dots, u_{n+1} .

△

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

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

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.

Solution. A necessary and sufficient condition in order for (3) to be satisfied for any polynomial $p \in \mathbf{P}(5)$ is that

$$\int_{-1}^1 x^d \, dx = \sum_{i=1}^n w_i x_i^d, \quad \text{for all } d \in \{0, 1, 2, 3, 4, 5\}.$$

This leads to the following system of equations

$$\begin{cases} 2 = w_1 + w_2 + w_3, \\ 0 = w_1 x_1 + w_2 x_2 + w_3 x_3, \\ \frac{2}{3} = w_1 x_1^2 + w_2 x_2^2 + w_3 x_3^2, \\ 0 = w_1 x_1^3 + w_2 x_2^3 + w_3 x_3^3, \\ \frac{2}{5} = w_1 x_1^4 + w_2 x_2^4 + w_3 x_3^4, \\ 0 = w_1 x_1^5 + w_2 x_2^5 + w_3 x_3^5. \end{cases}$$

Given the symmetry of the problem, it is reasonable to look for a solution of the form

$$(x_1, x_2, x_3, w_1, w_2, w_3) = (-x, 0, x, w_1, w_2, w_1),$$

where only 3 unknown parameters remain. For such a set of parameters, the second, fourth and sixth equations are satisfied, and the other three equations give

$$\begin{cases} 2 = 2w_1 + w_2, \\ \frac{2}{3} = 2w_1 x^2, \\ \frac{2}{5} = 2w_1 x^4. \end{cases}$$

Dividing the third equation by the second, we obtain $x^2 = 3/5$ and so $x = \pm\sqrt{3/5}$ (both values lead to the same integration rule in the end). It is then simple to deduce

that $w_1 = \frac{5}{9}$ and $w_2 = \frac{8}{9}$. We have thus derived the formula

$$\int_{-1}^1 u(x) \approx \frac{5}{9}u\left(-\sqrt{\frac{3}{5}}\right) + \frac{8}{9}u(0) + \frac{5}{9}u\left(\sqrt{\frac{3}{5}}\right).$$

△

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?

Solution. Clearly $L_0 = 1$. Then $L_1 = x + a_1$ and the requirement that $\langle L_1, L_0 \rangle = 0$ implies that $a_1 = 0$. We then use the ansatz $L_2 = x^2 + b_2x + a_2$ for L_2 . The requirement that $\langle L_2, L_1 \rangle$ leads to $b_2 = 0$, and then

$$\langle L_2, L_0 \rangle = \frac{2}{3} + 2a_2,$$

and so $L_2(x) = x^2 - \frac{1}{3}$. Finally, for L_3 , we use the ansatz $L_3 = x^3 + c_3x^2 + b_3x + a_3$. We calculate

$$\begin{aligned} \langle L_3, 1 \rangle &= \frac{2}{3}c_3 + 2a_3, \\ \langle L_3, x \rangle &= \frac{2}{5} + \frac{2}{3}b_3, \\ \langle L_3, x^2 \rangle &= \frac{2}{5}c_3 + \frac{2}{3}a_3. \end{aligned}$$

The second equation gives $b_3 = -\frac{3}{5}$, and the other two equations lead to $c_3 = a_3 = 0$. We conclude that $L_3(x) = x^3 - \frac{3}{5}x$. The roots of L_3 are given by $\left\{-\sqrt{\frac{3}{5}}, 0, \sqrt{\frac{3}{5}}\right\}$, and they coincide with the nodes of the Gauss–Legendre quadrature with 3 nodes. △

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 ℓ_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 .

Solution. Since (3) holds true for all $u \in \mathbf{P}(2n-1)$, it holds true in particular for the function $u = \ell_i \in \mathbf{P}(2n-1)$, which implies that

$$\int_{-1}^1 \ell_i(x) dx = \sum_{j=1}^n w_j \ell_i(x_j) = w_i.$$

Similarly, since (3) holds true also for $u \in \ell_i^2 \in \mathbf{P}(2n-1)$, we deduce that

$$\int_{-1}^1 (\ell_i(x))^2 dx = \sum_{j=1}^n w_j (\ell_i(x_j))^2 = w_i.$$

Since the left-hand side is positive, we deduce that $w_i > 0$. \triangle

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

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

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

$$\forall i \in \{1, \dots, n\}, \quad p(x_i) = u(x_i). \tag{5}$$

- 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

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

Solution. Assume that $p \in \mathbf{P}(2n-1)$ is such that (5) is satisfied. Then by (3) we deduce that

$$\int_{-1}^1 p(x) \, dx = \sum_{i=1}^n w_i p(x_i) = \sum_{i=1}^n w_i u(x_i) = \widehat{I}_n(u).$$

Consequently, we obtain that

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

This equation holds true in particular with p being the Hermite interpolation of u at the nodes x_1, \dots, x_n . Then, using the formula for the interpolation error, we obtain

$$u(x) - u(x) = \frac{u^{(2n)}(\xi(x))}{(2n)!} (x - x_1)^2 \dots (x - x_n)^2 = \frac{u^{(2n)}(\xi(x))}{(2n)!} (L_n(x))^2.$$

Indeed, as shown in class, L_n is a polynomial of degree n with single roots at x_1, \dots, x_n . Now we conclude by noting that

$$|I(u) - \widehat{I}_n(u)| = \left| \int_{-1}^1 u(x) - p(x) \, dx \right| \leq \int_{-1}^1 |u(x) - p(x)| \, dx \leq \int_{-1}^1 \frac{C_{2n}}{(2n)!} (L_n(x))^2 \, dx,$$

which concludes the exercise. \triangle

Question 4 (Vector and matrix norms, 10 marks). The 1-norm and the ∞ -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

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

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

1. (10 marks) $\|A\|_1$ is given by the maximum absolute column sum:

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|. \quad (6)$$

2. (Bonus +2) $\|A\|_\infty$ is given by the maximum absolute row sum:

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

Hint: In order to prove (6), 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 (6) by finding a vector \mathbf{x} with $\|\mathbf{x}\|_1 = 1$ such that

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

- Prove the direction \leq in (6) by showing that, for any $\mathbf{x} \in \mathbf{R}^n$ with $\|\mathbf{x}\|_1 = 1$,

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

Solution.

1. Let \mathbf{e}_j denote the column vector with a 1 at entry j and zero everywhere else. Notice

that $\|\mathbf{e}_j\|_1 = 1$ and

$$\|\mathbf{A}\mathbf{e}_{j_*}\|_1 = \sum_{i=1}^n |a_{ij_*}|,$$

and so $\|\mathbf{A}\|_1 \geq \sum_{i=1}^n |a_{ij_*}|$. It remains to prove that $\|\mathbf{A}\|_1 \leq \sum_{i=1}^n |a_{ij_*}|$. To this end, it is sufficient to show that $\|\mathbf{A}\mathbf{x}\|_1 \leq \sum_{i=1}^n |a_{ij_*}|$ for all $\mathbf{x} \in \mathbf{R}^n$ with $\|\mathbf{x}\|_1 = 1$. Take $\mathbf{x} \in \mathbf{R}^n$ with $\|\mathbf{x}\|_1 = 1$. We calculate that

$$\begin{aligned} \|\mathbf{A}\mathbf{x}\|_1 &= \sum_{i=1}^n \left| \sum_{j=1}^n a_{ij} x_j \right| \leq \sum_{i=1}^n \sum_{j=1}^n |a_{ij}| |x_j| \\ &= \sum_{j=1}^n \left(\sum_{i=1}^n |a_{ij}| \right) |x_j| \leq \sum_{j=1}^n \left(\sum_{i=1}^n |a_{ij_*}| \right) |x_j| \\ &= \left(\sum_{i=1}^n |a_{ij_*}| \right) \sum_{j=1}^n |x_j| = \left(\sum_{i=1}^n |a_{ij_*}| \right) \|\mathbf{x}\|_1 = \sum_{i=1}^n |a_{ij_*}|, \end{aligned}$$

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

2. Let i_* denote the index of a row (not necessarily unique) with maximum absolute sum, and let \mathbf{y} be a column vector with entry j equal to $\text{sign}(a_{i_*j})$. Then $\|\mathbf{y}\|_\infty = 1$ and

$$\|\mathbf{A}\mathbf{y}\|_\infty = \sum_{j=1}^n |a_{i_*j}|,$$

which implies that $\|\mathbf{A}\|_\infty \geq \sum_{j=1}^n |a_{i_*j}|$. It remains to prove that $\|\mathbf{A}\|_\infty \leq \sum_{j=1}^n |a_{i_*j}|$. To this end, take $\mathbf{x} \in \mathbf{R}^n$ with $\|\mathbf{x}\|_\infty = 1$. Then for all $i \in \{1, \dots, n\}$,

$$\begin{aligned} |(\mathbf{A}\mathbf{x})_i| &= \left| \sum_{j=1}^n a_{ij} x_j \right| \leq \sum_{j=1}^n |a_{ij}| |x_j| \leq \left(\sum_{j=1}^n |a_{ij}| \right) \max_{1 \leq j \leq n} |x_j| \\ &= \left(\sum_{j=1}^n |a_{ij}| \right) \|\mathbf{x}\|_\infty = \sum_{j=1}^n |a_{ij}| \leq \sum_{j=1}^n |a_{i_*j}|, \end{aligned}$$

which implies that $\|\mathbf{A}\|_\infty \leq \sum_{j=1}^n |a_{i_*j}|$.

△