

1. Mark each of the following statements T (True) or F (False).

Grading: 2 * (number correct) - 1.5 * (number incorrect) + 0 * (number blank).

- a. F Suppose A is an $n \times n$ nonsingular matrix with no zero elements. It takes fewer floating point operations to form A^{-1} and solve a linear system by multiplying $A^{-1}b$ than to factor the matrix $A = LU$ and use forward and backward substitution to solve the system.
- b. F The condition number of a matrix is always less than or equal to 1.
- c. T $\|A\|_1 = 5$ if

$$A = \begin{bmatrix} 2 & -1 \\ -3 & 1 \end{bmatrix}$$

- d. F As you increase the number n of data points, equally spaced in the interval $[0, 1]$, the error in the the composite Simpson's rule approximation to the integral of a smooth function is of order n to the power 4.
- e. T For any numerical integration routine with a nonrandom set of evaluation points, we can find a function whose integral is 1, but for which the routine reports an integral of zero.
- f. F When adding numbers, relative error bounds add.
- g. T The QR factorization usually gives a more accurate solution to a least squares problem than forming and solving the normal equations.
- h. F Backward error measures the distance between the true solution and the computed solution.
- i. F The software `ode15s` should not be used for stiff systems of differential equations.
- j. F In general, a cubic spline interpolant will oscillate more than a polynomial interpolant.

2a. (10) Suppose we have data points (t_i, y_i) , $i = 1, \dots, 20$, and we want to model the function $y(t)$ by a polynomial of degree 3 (i.e., a cubic polynomial) using a least squares fit. Give the dimensions $m \times n$ of the matrix that you would form in solving this least squares problem.

Answer: $m = 20$

$n = 4$

2b. (10) Compute $f[1, 2, 3]$ if $f(1) = 4$, $f(2) = 1$, and $f(3) = -5$.

Answer:

$$f[1, 2, 3] = \frac{\frac{4-1}{1-2} - \frac{-5-1}{3-2}}{1-3} = \frac{-3+6}{-2} = -\frac{3}{2}$$

3. Suppose A is a matrix of size 100×50 ,

B is a matrix of size 50×50 , and

x is a vector of size 50×1 .

3a. (10) How many floating point multiplications does it take to form $A * (B * x)$?

Answer: $(50*50)$ to form $B * x$, and then $100 * 50$ to form A times this.

3b. (10) How many floating point multiplications does it take to form $(A * B) * x$?

Answer: $(100*50*50)$ to form $A * B$ and then $100 * 50$ to form this times x . (This is 34 times the work of 3a.)

4a. (10) Suppose we have 20 data points (x_i, y_i) , with $x_1 < x_2 < \dots < x_n$. Give a formula that approximates

$$Q = \int_{x_1}^{x_n} y(x) dx .$$

Answer: Use Trapezoidal, since we don't know whether the spacing is equal.

$$Q \approx \sum_{i=1}^{19} \frac{x_{i+1} - x_i}{2} (y_{i+1} + y_i) .$$

4b. (10) Suppose that $f(x)$ is given as a Matlab m-file. How could you use `ode45` to approximate

$$\int_1^8 f(x) dx?$$

Answer: Solve the differential equation

$$y' = f(x), \quad y(1) = 0.$$

The integral will be $y(8)$, so we need the last value from the y vector in `[x,y] = ode45('f',[1,8],0)`.

5. (20) A scientist has solved a linear system on a machine with “machine epsilon” equal to 10^{-16} . The matrix was of size 100×100 , the condition number was 10^3 , the residual was $\|b - Ax\| = 10^{-12}$, and the algorithm was Gaussian elimination with partial pivoting (i.e., the Matlab backslash operator). The data has bounds $\|b\| = 70$, $\|x_{comp}\| = 10$, and $\|A\| = 30$. The values b_i come from measured data, and the difference between b_i and its true value d_i is bounded by $|b_i - d_i| \leq 10^{-5}$, $i = 1, \dots, 100$. Give the scientist forward and backward error bounds on the results she computed compared to the problem she really wanted to solve: $Ay = d$. Explain to her what these bounds mean about the quality of her results.

Answer:

1. Backward error: The residual $d - Ax$ is bounded by $100 * 10^{-5}$ (since for the norm, we need to add 10^{-5} 100 times), so she has solved the problem $Ay = d + e$, where $\|e\| < 10^{-3}$.
2. Forward error:

$$\frac{\|x - x_{true}\|}{\|x_{true}\|} \leq \kappa \frac{\|e\|}{\|d\|} \leq 10^3 \frac{10^{-3}}{70} = \frac{1}{70}.$$

This is how close the computed answer is to the true answer.

6. (25) Suppose we have an increasing function $y = f(x)$, and suppose that we have stored two vectors \mathbf{x} and \mathbf{y} in Matlab, with $\mathbf{y}(i) = f(\mathbf{x}(i))$. Suppose further that $\mathbf{x}(1) = -4$, $\mathbf{y}(1) = -2$, $\mathbf{x}(5) = 6$, $\mathbf{y}(5) = 10$. Use the Matlab functions for cubic spline interpolation to determine a point x for which $f(x) \approx 0$.

Note: Up to full credit for “inverse interpolation.” Up to 15 points for interpolation plus root finding.

Note: f is increasing if $f(x) > f(z)$ whenever $x > z$.

Answer: The approximate zero is $z = \text{spline}(\mathbf{y}, \mathbf{x}, 0)$.

7. Refer to the listing of `fzero` that is attached after this page.

7a. (15) Mark changes on the listing to add another output variable `count`, which is a count of the number of function evaluations that were performed. (Don't forget to change the documentation.)

Answer: Initialize `count` to zero, add one to it after each call to `feval`, and add documentation explaining the purpose of `count` and documenting who changed the program and when.

7b. (10) Why did the author use `p` and `q` rather than just defining `d` directly?

Answer: If s is too close to 1, for example, then p/q can overflow, so we want to avoid making the division until we check (in the line following “Is interpolated point acceptable”) that it is safe.