# MS/PhD Qualifying exam: Numerical Analysis August 22, 2014

**Closed book/closed notes.**
**All questions are equally weighted.**
**Show all working.**

## Part 1 (MATH:5800/22M:170)

1. *Floating point arithmetic.* The standard formula for computing the roots of a quadratic $ax^2 + bx + c = 0$:

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

   is known to have problems computing the smaller root numerically if $b^2$ is *much* larger than $|ac|$. Explain the cause of the cause of this problem; propose a method for avoiding this problem.

2. *Solution of nonlinear equations.* Carry out two steps of the secant method for solving $x - \cos x = 0$ starting with $x_0 = 0$ and $x_1 = 1$. What rate of convergence is expected, and under what conditions is this rate of convergence obtained?

3. *Interpolation and approximation.* Using equally spaced interpolation points is known to result in Runge's phenomenon for the function $f(x) = 1/(1 + x^2)$ interpolated over $[-5, +5]$. What is this phenomenon? Can the use of a different set of interpolation points prevent this phenomenon? If so, explain how?

4. *Numerical integration.* Use Simpson's method with five function evaluations to obtain an estimate of $\int_0^1 e^x/(1 + x)\, dx$. What is the asymptotic order of the error of composite Simpson's method with $2n + 1$ function evaluations? Give an example of a method that has an asymptotically faster rate of convergence than Simpson's method as the number of function evaluations goes to infinity.

# Part 2 (MATH:5810/22M:171)

1. *Multistep methods.* Consider the general multistep method

$$y_{n+1} = \sum_{j=0}^{p} a_j y_{n-j} + h \sum_{j=-1}^{p} b_j f(t_{n-j}, y_{n-j}).$$

   In order to prove convergence of a particular order for this method we need two basic conditions: a stability condition, and a consistency condition. Give these conditions. Use them to determine if, and with what order, the leap-frog method converges:

$$y_{n+1} = y_{n-1} + 2h f(t_n, y_n).$$

2. *Runge–Kutta methods.* Show that Heun's method

$$\begin{aligned} \mathbf{z}_{n+1} &= \mathbf{y}_n + h\, \mathbf{f}(t_n, \mathbf{y}_n), \\ \mathbf{y}_{n+1} &= \mathbf{y}_n + \frac{1}{2}h\left[ \mathbf{f}(t_n, \mathbf{y}_n) + \mathbf{f}(t_{n+1}, \mathbf{z}_{n+1}) \right] \end{aligned}$$

   has a local truncation error of $\mathcal{O}(h^3)$. What is its asymptotic global truncation error in the form $\mathcal{O}(h^m)$?

3. *LU factorization and linear systems.* The perturbation theorem for linear systems states that if $Ax = b$, $(A+E)\widehat{x} = b + d$, and $\left\| A^{-1} \right\| \|E\| < 1$, then

$$\frac{\|\widehat{x} - x\|}{\|x\|} \leq \frac{\kappa(A)}{1 - \kappa(A)(\|E\| / \|A\|)} \left[ \frac{\|E\|}{\|A\|} + \frac{\|d\|}{\|b\|} \right]$$

   where $\kappa(A) = \left\| A^{-1} \right\| \|A\|$ is the condition number. Using this, how many digits of accuracy are expected in the computed solution $\widehat{x}$ given that the matrix $A$ and right-hand side $b$ are known to 5 digits, but $\kappa(A) \approx 10^3$? The backward error theory for *LU* factorization by Wilkinson shows that the computed solution $\widehat{x}$ of a system $Ax = b$ exactly satisfies $(A+E)\widehat{x} = b$ where $\|E\|_\infty \leq 3\mathbf{u} \left( \|A\|_\infty + \left\| \widehat{L} \right\|_\infty \left\| \widehat{U} \right\|_\infty \right)$ where $\widehat{L}$ and $\widehat{U}$ are the computed $L$ and $U$ factors in the *LU* factorization. If $\left\| \widehat{L} \right\|_\infty \left\| \widehat{U} \right\|_\infty / \|A\|_\infty$ is modest (say $\approx 10$), give an estimate for the relative error $\|\widehat{x} - x\|_\infty / \|x\|_\infty$ in terms of $\kappa(A)$ in the $\infty$-norm.

4. *Least squares and QR factorization.* What is a QR factorization of a matrix? Describe two different ways of computing the QR factorization of an $m \times n$ matrix. Explain how to use a QR factorization of a matrix to solve a least squares problem $\min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|_2$.

# NUMERICAL ANALYSIS EXAM

<u>Directions:</u> Answer all the 6 problems.

1. For a linearly convergent iteration $x_{n+1} = g(x_n)$, $g$ being continuously differentiable, we have $x_{20} = 1.3254943$, $x_{21} = 1.3534339$, $x_{22} = 1.3708962$. Show how to estimate (you do not need to compute the numbers)

   (a) the fixed-point $\alpha$ of the function $g$;

   (b) the rate of linear convergence;

   (c) the error $\alpha - x_{22}$.

   *Hint*: From the assumption, there is a constant $\lambda$ such that for $n$ large, $(x_{n+1}-\alpha)/(x_n-\alpha) \approx \lambda$.

2. Let $f \in C([0,1])$ be given and let $0 = x_0 < x_1 < \cdots < x_{N-1} < x_N = 1$ be a partition of the interval $[0,1]$. Denote by $s$ the piecewise linear interpolant of $f$ corresponding to the partition; i.e., $s(x)$ is a linear function on each subinterval $[x_{n-1}, x_n]$, $n = 1, \ldots, N$, and $s(x_n) = f(x_n)$, $n = 0, 1, \ldots, N$.

   (a) Give a formula for $s$ on each subinterval.

   (b) Assuming $f \in C^2([0,1])$, bound the error $f(x) - s(x)$.

3. (a) Find the constant $c$ that minimizes $\max_{0 \le x \le 1} |e^x - c|$.

   (b) Find the constant $c$ that minimizes $\displaystyle\int_0^1 |e^x - c|^2 dx$.

   (c) Find an equation for the constant $c$ that minimizes $\displaystyle\int_0^1 |e^x - c| \, dx$.

4. Consider solving the initial value problem $y' = f(x, y)$ for $0 \le x \le 1$, $y(0) = Y_0$, $f$ being a smooth function. Let $0 = x_0 < x_1 < \cdots < x_N = 1$ be a uniform partition of the interval $[0,1]$ and denote $h$ the step size. For a constant parameter $\theta \in [0,1]$, introduce the following generalized mid-point method

$$y_{n+1} = y_n + h \left[ (1 - \theta) f(x_n, y_n) + \theta f(x_{n+1}, y_{n+1}) \right].$$

   It is known that for $h$ small enough, this relation defines a unique value $y_{n+1}$.

   (a) Determine the order of the method.

   (b) Show that the method is absolutely stable when $\theta \in [1/2, 1]$.

5. What is the Cholesky factorization? Find the Cholesky factorization of the matrix

$$A = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 5 & 5 \\ 1 & 5 & 14 \end{pmatrix}.$$

6. In iteratively solving the linear system $A\boldsymbol{x} = \boldsymbol{b}$ ($\det A \neq 0$), we write $A = P - N$ with $P$ nonsingular, and generate a sequence $\{\boldsymbol{x}^{(k)}\}$ by the formula

$$P\boldsymbol{x}^{(k+1)} = \boldsymbol{b} + N\boldsymbol{x}^{(k)},$$

starting with some initial guess $\boldsymbol{x}^{(0)}$. Denote the residual $\boldsymbol{r}^{(k)} = \boldsymbol{b} - A\boldsymbol{x}^{(k)}$.
(a) Show that the iteration formula can be equivalently expressed as

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} + P^{-1}\boldsymbol{r}^{(k)}.$$

(b) Let $\alpha > 0$ be a constant. Define the stationary Richardson method by the formula

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} + \alpha P^{-1}\boldsymbol{r}^{(k)}.$$

Show that the method converges if and only if $\alpha|\lambda|^2 < 2\operatorname{Re}\lambda$ for any eigenvalue $\lambda$ of $P^{-1}A$.

# PhD/MS Qualifying Exam        Numerical Analysis

**Department of Mathematics**                    **University of Iowa**

**116 MH**                        **2:30–5:30 pm, January 23, 2008**

1. Here is some MATLAB$^{TM}$ code to evaluate $e^x$.

   ```
   function y = myexp(x)
   % Computes y = e^x using Taylor series
   u = 2.2e-16; % unit roundoff
   term = 1;   y = 0;   k = 0;
   while abs(term) > u
       y = y + term;   term = x*term/(k+1);
       k = k+1;
   end
   ```

   This code gives accurate values for small $x$. But when $x = -20$ is used as the input, the computed value is $\approx 5.62 \times 10^{-9}$, while `exp(-20)` gives $\approx 2.06 \times 10^{-9}$; there is not even one correct digit. Investigation shows that during the computation of `myexp(-20)`, the value of `term` becomes as large as $\approx 4.3 \times 10^7$. Use your knowledge of floating point arithmetic (in double precision) to explain why the error is so large. This appears to be a problem for large negative values of $x$. How can you modify this code to give accurate values for all $x$ (ignoring the problems of over- and under-flow)?

2. Write out formulas or pseudo-code for the Newton and secant methods for solving $f(x) = 0$. List all conditions needed in order for these methods to converge. Give the expected rates of convergence of these methods under the conditions you have stated. Explain your terms. Perform three steps of Newton's method to find the positive root of $x^3 - 3x - 3 = 0$ using a starting value of $x_0 = 2$.

3. Given a function $f(x)$ on an interval $[a, b]$ and points $a \leq x_0 < x_1 < x_2 < \cdots < x_n \leq b$, give a method to construct a polynomial $p(x)$ of degree $\leq n$ where $p$ interpolates $f$ at $x_0, x_1, \ldots, x_n$. (You do not have to give pseudo-code, but it should be a clear and complete description.) Give a formula for estimating the interpolation error $f(x) - p(x)$.

1

What is Chebyshev interpolation? From the formula for the interpolation error, explain how Chebyshev interpolation relates to minimax approximation.

4. Consider solving the initial value problem $y' = f(x, y)$ for $0 \leq x \leq 1$, $y(0) = Y_0$, $f$ being a smooth function. Let $0 = x_0 < x_1 < \cdots < x_N = 1$ be a uniform partition of the interval $[0, 1]$ and denote $h$ the step size. Consider a method of the form

$$y_{n+1} = \alpha\, y_n + \beta\, y_{n-1} + h\, \gamma\, f(x_{n-1}, y_{n-1}), \quad n \geq 1;$$
$$y_0 = Y_0, \quad y_1 = Y_0 + h\, f(x_0, Y_0).$$

Choose the constants $\alpha$, $\beta$, and $\gamma$ so that the order of the method is as high as possible. Determine whether the resulting method is convergent.

5. Define the Gauss-Jacobi method and the Gauss-Seidel method for solving the linear system $A\mathbf{x} = \mathbf{b}$, where $\mathbf{b} \in \mathbb{R}^N$ is given and

$$A = \begin{pmatrix} 2 & -1 & & & \\ -1 & 3 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 3 & -1 \\ & & & -1 & 2 \end{pmatrix} \in \mathbb{R}^{N \times N}.$$

For each method, determine a number of iterations to reduce the $\|\cdot\|_\infty$ norm of the error by a factor of $0.01$.

6. What is the singular value decomposition (SVD) of a general rectangular matrix $A \in \mathbb{C}^{m \times n}$? For a $3 \times 3$ matrix $B$, compute $\|B\|_2$ and $\mathrm{Cond}_2(B) = \|B\|_2\|B^{-1}\|_2$ by taking advantage of the following SVD of the matrix $B$ (four significant digits after the decimal point are displayed):

$$\begin{pmatrix} -0.3150 & -0.4193 & -0.8515 \\ -0.6449 & -0.5637 & 0.5161 \\ -0.6963 & 0.7117 & -0.0928 \end{pmatrix} \begin{pmatrix} 10.2377 & 0 & 0 \\ 0 & 4.4830 & 0 \\ 0 & 0 & 0.3050 \end{pmatrix}$$
$$\times \begin{pmatrix} -0.5598 & 0.3230 & 0.7630 \\ -0.7216 & 0.2625 & -0.6406 \\ -0.4073 & -0.9092 & 0.0861 \end{pmatrix}.$$

# MS Exam on Numerical Analysis

Directions: Answer all the 6 problems.

1. Give the Newton method for solving the equation

$$e^x + 4\,e^{-x} - 4 = 0,$$

   and discuss the convergence order of the method.

2. Let $f(x) = \sin(\pi x)$. Determine a function $p(x)$ such that $p(x)$ is a polynomial on $[0, 0.5]$ and $[0.5, 1]$, and satisfies the conditions

$$p(x) = f(x), \ p'(x) = f'(x), \quad \text{for } x = 0, 0.5, 1.$$

3. (a) Find $A_0$, $A_1$ and $A_2$ such that the integration rule

$$I(f) = \int_{-h}^{h} f(x)\, dx \approx A_0 f\left(-h/2\right) + A_1 f(0) + A_2 f\left(h/2\right)$$

is exact for polynomials of degree $\leq 2$.

(b) Show that the rule constructed in (a) is in fact exact for polynomials of degree $\leq 3$.

(c) For the constructed rule, it can be proved that

$$I(f) - \left(A_0 f\left(-h/2\right) + A_1 f(0) + A_2 f\left(h/2\right)\right) = c_0 f^{(4)}(\eta)\, h^5, \quad \eta \in [-h, h]$$

where $c_0$ is a constant independent of $f$. Find the constant $c_0$.

4. Suppose that $B \approx A^{-1}$. Starting with an initial guess $x_0$, consider the following residual correction method:

$$\text{for } k = 0, 1, 2, \ldots$$
$$r_k \leftarrow Ax_k - b;$$
$$x_{k+1} \leftarrow x_k - Br_k.$$

Show that this will converge (using exact arithmetic) so that $x_k \to x$, with $x$ the exact solution, provided $\|I - BA\| < 1$.

5. What is the QR factorization of a matrix? Explain how a QR factorization of a matrix can be computed using any of the following three methods: (i) Gram–Schmidt orthgonalization, (ii) Givens' rotations, or (iii) Householder reflectors.

An overdetermined linear system $Ax = b$ where $A$ is $m \times n$ with $m > n$ usually cannot be solved for $x$. Instead we can ask to minimize $\|Ax - b\|_2$ over all $x$. Show how the solution to this minimization problem can be computed using the QR factorization of $A$.

6. Consider the following two methods for numerically solving an initial value problem for the ODE $dx/dt = f(t, x)$:

$$x_{n+1} = x_{n-1} + 2\,h\,f(t_n, x_n) \qquad \text{(leap-frog method)}$$

$$x_{n+1} = x_n + (h/2)[f(t_n, x_n) + f(t_{n+1}, x_{n+1})] \qquad \text{(implicit mid-point method)}$$

where $h$ is the step size, and $t_k = t_0 + k\,h$. For the test equation $dx/dt = \lambda x$, show that the leap-frog method is only stable for $\lambda h = 0$, while the implicit mid-point method is stable for all $\lambda h < 0$.

# MS/PhD Qualifying exam: Numerical Analysis
# August 20, 2020

**Closed book/closed notes.**
**All questions are equally weighted.**
**Show all working.**
**Bring a calculator.**
**No communication devices.**

## Part 1 MATH:5800

1. *Floating point arithmetic.* In Matlab, the formulas

$$\cos\theta = \frac{x}{\sqrt{x^2+y^2}}, \qquad \text{and}$$
$$\sin\theta = \frac{y}{\sqrt{x^2+y^2}}$$

give $\cos\theta = \sin\theta = 0$ for $x = y = 10^{+200}$. Why does this occur? Can you re-arrange the formulas to give equivalent formulas in exact arithmetic, but give accurate values for $\cos\theta$ and $\sin\theta$.

2. *Solution of nonlinear equations.* Carry out two steps of the secant method for solving $x\,e^x = 3$ starting with $x_0 = 0$ and $x_1 = 1$. What rate of convergence is expected, and under what conditions is this rate of convergence obtained?

3. *Interpolation and approximation.* Using equally spaced interpolation points is known to result in Runge's phenomenon for the function $f(x) = 1/(1+x^2)$ interpolated over $[-5, +5]$. What is this phenomenon? Can the use of a different set of interpolation points prevent this phenomenon? If so, explain how?

4. *Numerical integration.* Use Simpson's method with five function evaluations to obtain an estimate of $\int_0^1 e^x/(1+x)\,dx$. What is the asymptotic order of the error of composite Simpson's method with $2n+1$ function evaluations? Give an example of a method that has an

asymptotically faster rate of convergence than Simpson's method as the number of function evaluations goes to infinity.

# Part 2 MATH:5810

1. *Multistep methods.* Consider the general multistep method

$$y_{n+1} = \sum_{j=0}^{p} a_j y_{n-j} + h \sum_{j=-1}^{p} b_j f(t_{n-j}, y_{n-j})$$
$$= \sum_{j=0}^{p} a_j y_{n-j} + h \sum_{j=-1}^{p} b_j y'_{n-j}.$$

In order to prove convergence of a particular order for this method we need two basic conditions: a stability condition, and a consistency condition. [**Hint:** The consistency conditions come from Taylor series expansions of $y_{n-j} = y(t_{n-j})$. Give these conditions. Use them to determine if, and with what order, the trapezoidal method converges:

$$y_{n+1} = y_{n-1} + \frac{h}{2} \left[ f(t_n, y_n) + f(t_{n+1}, y_{n+1}) \right].$$

*Correction:* $y_{n-1}$ *on the right should be* $y_n$.

$$y_{n+1} = y_n + \frac{h}{2} \left[ f(t_n, y_n) + f(t_{n+1}, y_{n+1}) \right].$$

*End of correction.*

2. *Runge–Kutta methods.* Show that Heun's method

$$\mathbf{z}_{n+1} = \mathbf{y}_n + h\,\mathbf{f}(t_n, \mathbf{y}_n),$$
$$\mathbf{y}_{n+1} = \mathbf{y}_n + \frac{1}{2}h \left[ \mathbf{f}(t_n, \mathbf{y}_n) + \mathbf{f}(t_{n+1}, \mathbf{z}_{n+1}) \right]$$

has a local truncation error of $\mathcal{O}(h^3)$. What is its asymptotic global truncation error in the form $\mathcal{O}(h^m)$?

3. *LU factorization and linear systems.* The perturbation theorem for linear systems states that if $Ax = b$, $(A + E)\hat{x} = b + d$, and $\|A^{-1}\| \, \|E\| < 1$, then

$$\frac{\|\hat{x} - x\|}{\|x\|} \leq \frac{\kappa(A)}{1 - \kappa(A)(\|E\| / \|A\|)} \left[ \frac{\|E\|}{\|A\|} + \frac{\|d\|}{\|b\|} \right]$$

where $\kappa(A) = \|A^{-1}\| \, \|A\|$ is the condition number. Using this, how many digits of accuracy are expected in the computed solution $\hat{x}$ given that the matrix $A$ and right-hand side $b$ are known to 5 digits, but $\kappa(A) \approx 10^3$? The backward error theory for $LU$ factorization by Wilkinson shows that the computed solution $\hat{x}$ of a system $Ax = b$ exactly satisfies $(A + E)\hat{x} = b$ where $\|E\|_\infty \leq 3\mathbf{u} \left( \|A\|_\infty + \left\| \hat{L} \right\|_\infty \left\| \hat{U} \right\|_\infty \right)$ where $\hat{L}$ and $\hat{U}$ are the computed $L$ and $U$ factors in the $LU$ factorization. If $\left\| \hat{L} \right\|_\infty \left\| \hat{U} \right\|_\infty / \|A\|_\infty$ is modest (say $\approx 10$), give an estimate for the relative error $\|\hat{x} - x\|_\infty / \|x\|_\infty$ in terms of $\kappa(A)$ in the $\infty$-norm.

4. *QR algorithm.* Give the QR algorithm (with *or* without shifting). Show that if the original matrix $A$ is symmetric, then every iterate of the QR factorization is symmetric. Explain how shifting is used to improve the rate of convergence of the QR algorithm.

# MS/PhD Qualifying exam
# Numerical Analysis II (22M:171)
# August 22, 2014

**Closed book/closed notes.**
**All questions are equally weighted.**
**Show all working.**

1. *Multistep methods.* Consider the general multistep method

$$y_{n+1} = \sum_{j=0}^{p} a_j y_{n-j} + h \sum_{j=-1}^{p} b_j f(t_{n-j}, y_{n-j}).$$

In order to prove convergence of a particular order for this method we need two basic conditions: a stability condition, and a consistency condition. Give these conditions. Use them to determine if, and with what order, the leap-frog method converges:

$$y_{n+1} = y_{n-1} + 2h f(t_n, y_n).$$

2. *Runge–Kutta methods.* The implicit trapezoidal rule is

$$y_{n+1} = y_n + \frac{1}{2} h \left[ f(t_n, y_n) + f(t_{n+1}, y_{n+1}) \right].$$

Describe what is meant by the stability region of a Runge–Kutta method. What is the stability region for the implicit trapezoidal method? Rigorously justify your answer.

3. *LU factorization and linear systems.* The perturbation theorem for linear systems states that if $Ax = b$, $(A+E)\widehat{x} = b+d$, and $\left\| A^{-1} \right\| \|E\| < 1$, then

$$\frac{\|\widehat{x} - x\|}{\|x\|} \leq \frac{\kappa(A)}{1 - \kappa(A)(\|E\| / \|A\|)} \left[ \frac{\|E\|}{\|A\|} + \frac{\|d\|}{\|b\|} \right]$$

where $\kappa(A) = \left\| A^{-1} \right\| \|A\|$ is the condition number. Using this, how many digits of accuracy are expected in the computed solution $\widehat{x}$ given that the matrix $A$ and right-hand side $b$ are known to 5 digits, but $\kappa(A) \approx 10^3$?

The backward error theory for $LU$ factorization by Wilkinson shows that the computed solution $\widehat{x}$ of a system $Ax = b$ exactly satisfies $(A + E)\widehat{x} = b$ where $\|E\|_\infty \le 3\mathbf{u}\left(\|A\|_\infty + \left\|\widehat{L}\right\|_\infty \left\|\widehat{U}\right\|_\infty\right)$ where $\widehat{L}$ and $\widehat{U}$ are the computed $L$ and $U$ factors in the $LU$ factorization. If $\left\|\widehat{L}\right\|_\infty \left\|\widehat{U}\right\|_\infty / \|A\|_\infty$ is modest (say $\approx 10$), give an estimate for the relative error $\|\widehat{x} - x\|_\infty / \|x\|_\infty$ in terms of $\kappa(A)$ in the $\infty$-norm.

4. *Least squares, Cholesky and QR factorization.* The normal equations for solving a least squares problem $\min_x \|Ax - b\|_2$ are $A^T A x = A^T b$. Describe precisely what the Cholesky and $QR$ factorizations are. Show how to solve the least squares problem using either the $QR$ factorization applied to the original system, or Cholesky factorization applied to the normal equations for the least squares problem. Also show that the $R$ matrix in the $QR$ factorization of $A$ is one of the factors in a Cholesky factorization of $A^T A$.