

## 1.1 Basic concepts

Absolute and relative error:

Absolute =  $AF(a) = \text{approx.}(a) - \text{true}(a)$

Relative =  $RF(a) = AF(a) / \text{true}(a)$

precision: the number of digits

accuracy: the number of correct significant digits

Truncation error: Difference between true result and result given by exact arithmetic.

Rounding error: Difference between true result and result by rounded arithmetic

## 1.2 Floating-point numbers

Floating point is characterized by:

$\beta$ : Base or radix

$p$ : precision

$[L, U]$ : exponent range

$$\Rightarrow \text{float} = \pm \left( d_0 + \frac{d_1}{\beta} + \frac{d_2}{\beta^2} + \dots + \frac{d_{p-1}}{\beta^{p-1}} \right) \beta^E \quad \text{with: } 0 \leq d_i < \beta-1 \quad L \leq E \leq U$$

normalized if:  $d_0 = 1$

(Binary system has  $\beta=2$ )

## 1.3 Properties

float is finite and discrete.

smallest positive normalized number (underflow level):

$\beta^L$

largest (overflow):  $\beta^{U+1}(1-\beta^{-P})$

Real number representable in given floating point

system; Machine numbers.

$\Rightarrow$  otherwise round to nearest float; rounding error  
accuracy of floating-point system is machine precision.

## 1.4 Good practice for computer arithmetic

cancellation:

- avoid subtracting 2 almost identical numbers

Addition:

- avoid adding small and large numbers
- Perform a sequence of additions ordered from the smallest to the largest.

## 2. Linear systems and notation

Linear system:  $Ax + b$

$$= \begin{cases} a_{11}x_1 + a_{12}x_2 \dots = b_1 \\ a_{21}x_1 + a_{22}x_2 \dots = b_2 \\ \vdots \\ a_{m1}x_1 + \dots + a_{mn}x_n = b_m \end{cases}$$

An  $n \times n$  matrix is nonsingular if:

- $A$  has inverse  $A^{-1} \Rightarrow AA^{-1} = I$
- $\det(A) \neq 0$
- $\text{rank}(A) = n$
- $x \neq 0 \Rightarrow Ax \neq 0$

## 2.2. Solving linear systems

To solve we transform  $Ax = b$  in another one whose solution is easy to compute:

We can multiply by a nonsingular matrix without effecting the solution.

Easy to calculate?  $\Rightarrow$  triangular linear system.

A matrix  $L$  is lower triangular if  $l_{ij} = 0$  if  $i < j$

$\Rightarrow$  solved by forward substitution:

$$Lx = b \quad \stackrel{\text{forward}}{\Rightarrow} \quad x_1 = b_1 / l_{11}; x_i = (b_i - \sum_{j=1}^{i-1} l_{ij}x_j) / l_{ii}$$

A matrix  $U$  is upper triangular if  $u_{ij} = 0$  if  $i > j$

$\Rightarrow$  solved by backward substitution:

$$Ux = b \Rightarrow x_m = b_m / u_{mm}, x_i = (b_i - \sum_{j=1}^m u_{ij}x_j) / u_{ii}$$

## Elementary elimination matrices.

A matrix (nonsingular) that transforms a system into a triangular linear system.

Gaussian elimination = LU factorization / decomposition.

- ⇒ • Use elementary elimination matrices to that the system transforms into upper triangular.
- $M = M_{n-1} \cdots M_1 \Rightarrow MAx = Mb$  upper triangular.  
 $L = M^{-1}$  lower triangular.
- ⇒  $U = MA$ ,  $L = M^{-1}$
- ⇒  $LUx = b$  solved by first:  $Ly = b$   
 then:  $Ux = y \Rightarrow \text{X}$

## Partial pivoting

### Gauss-Jordan elimination

- ⇒ variation of gaussian elimination
- ⇒ eliminates entries below and above diagonal entries
- ⇒ 50% more computationally expensive than Gauss-el.

## 9.3 special types of linear systems

- Symmetric  $A = A^T$
- positive definite  $x^T Ax > 0$  for all  $x \neq 0$
- Banded:  $a_{ij} = 0$  for all  $\|i-j\| \geq \beta$ ,  $\beta = \text{bandwidth of } A$
- sparse: most entries are zero.

Symmetric positive definite systems:

### Cholesky factorization

if  $A$  is symmetric and positive definite:  $V = L^T \Rightarrow A = LL^T$

$\Rightarrow$  Cholesky factorization: compute  $LL^T \Rightarrow$  easily computable.

### Computational complexity

- LU factorization takes about  $\frac{m^3}{3}$  floating point operations
- A complete matrix about  $m^3$
- Solving LU about  $m^2$

## 2.4 Sensitivity and conditioning

### Vector norms

$$p\text{-norms} \Rightarrow \|x\|_p = \left( \sum_{i=1}^m \|x_i\|^p \right)^{1/p}$$

$$\begin{aligned} \Rightarrow & \cdot 1\text{-norm: } \|x\|_1 = \sum_{i=1}^m \|x_i\| \\ & \cdot 2\text{-norm: } \|x\|_2 = \left( \sum_{i=1}^m \|x_i\|^2 \right)^{1/2} \\ & \cdot \infty\text{-norm: } \|x\|_\infty = \max_{1 \leq i \leq m} \|x_i\| \end{aligned}$$

In general:

- $\|x\|_1 \geq \|x\|_2 \geq \|x\|_\infty$
- $\|x\|_1 \leq \sqrt{m} \|x\|_2$
- $\|x\|_2 \leq \sqrt{m} \|x\|_\infty$

For all  $p$ -norms:

- $\|x\| > 0$  if  $x \neq 0$
- $\|\alpha x\| = \|\alpha\| \cdot \|x\|$
- $\|x+y\| \leq \|x\| + \|y\|$

### Matrix norms:

$$\|A\| = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|}$$

$$\|A\|_1 = \max_j \sum_{i=1}^m \|a_{ij}\| \quad \& \quad \|A\|_\infty = \max_i \sum_{j=1}^m \|a_{ij}\|$$

Properties:

- $\|A\| \geq 0$  if  $A \geq 0$
- $\|\alpha A\| = \|\alpha\| \cdot \|A\|$
- $\|A+B\| \leq \|A\| + \|B\|$
- $\|AB\| \leq \|A\| \cdot \|B\|$
- $\|Ax\| \leq \|A\| \cdot \|x\|$

Matrix condition number.

$$\text{cond}(A) = \|A\| \cdot \|A^{-1}\|$$

Error estimation.

$$Ax = b \quad \text{and} \quad Ax' = b + \Delta b$$

$$A\Delta x = \Delta b$$

$$\Rightarrow \frac{\|\Delta x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\Delta b\|}{\|b\|}$$

Residual:

Residual  $r$  of approx. solution  $x'$  of  $Ax = b$

$$\Rightarrow r = b - Ax'$$

In theory for nonsingular  $A$ :  $r = 0$  and  $\Delta x = \|x - x'\| = 0$

$\Rightarrow$  In practice not!

$$\Rightarrow \frac{\|\Delta x\|}{\|x'\|} \leq \text{cond}(A) \frac{\|r\|}{\|A\| \cdot \|x\|}$$

## 25 Software

scipy.linalg.solve, linalg.Cu, linalg.chol

## 2.1 Introduction

An overdetermined problem  $Ax = b$  with  $A: m \times n$  with  $m > n$

minimize the residual:  $r = b - Ax$

## 2.2 Normal equations

To minimize  $r = b - Ax$

We have:  $A^T A x = A^T b$  if  $\text{rank}(A) = n$

## 2.3 Problem transformations

Transform it into a square linear system.

### QR orthogonal transformations

$Q$  is orthogonal if  $Q^T Q = I$ , it preserves the Euclidean norm ( $2$ -norm) of a vector

### Triangular least squares problems.

$$\begin{bmatrix} R \\ 0 \end{bmatrix} \hat{x} \cong [c_1 c_2] \Rightarrow \|r\|_2^2 = \|c_1 - Rx\|_2^2 + \|c_2\|_2^2$$

$$\text{if } Rx = c_1 \Rightarrow \|r\|_2^2 = \|c_2\|_2^2$$

### QR-factorization

$$A = Q \begin{bmatrix} R \\ 0 \end{bmatrix} ; Q^T b = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}$$

$$\Rightarrow \|r\|_2^2 = \|c_1 - Rx\|_2^2 + \|c_2\|_2^2$$

$$\text{if } Rx = c_1 \Rightarrow \|r\|_2^2 = \|c_2\|_2^2$$

## Housholder transformations

$$H = I - 2 \frac{vv^T}{v^Tv}$$

$$H = H^T = HT$$

The product of Housholder transformations in itself an orthogonal matrix:

$$Q^T = H_m \dots H_1 \quad \text{or} \quad Q = H_m \dots H_1^T$$

## 2.4 Rank deficiency

If  $\text{rank}(A) < m \Rightarrow$  will perform QR factorization

but Upper triangular matrix will be singular.

$\Rightarrow$  multiple x vectors have the same minimal norm

## 2.5 SVD and Diagonalization

Diagonal Linear Least Squares eigenvectors value.

SVD of A:  $A = U\Sigma V^T$

$$\Sigma: \sigma_{ij} := \begin{cases} 0 & i \neq j \\ \sigma_i & i = j \end{cases}$$

$\sigma_i$  are the singular values of A

columns  $U_i$  of U and  $v_i$  of V are singular vectors

useful for ill-conditioned or nearly rank-deficient problems.

Other applications:

Euclidean matrix norm:  $\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \sigma_{\max}$

Euclidean condition number:  $\text{cond}_2(A) = \frac{\sigma_{\max}}{\sigma_{\min}}$

Rank determination: = number of singular values.

## Bseudo inverse

- Of a scalar:  $\sigma \Leftrightarrow \frac{1}{\sigma}$
- of a Diagonal matrix: Transpose and ( $\rho$ )<sup>-1</sup> of each entry

$\Rightarrow$  A pseudoinverse:

$$\tilde{A}^+ = V \Sigma^+ U^T$$

if  $A$  is square and nonsingular:  $\tilde{A}^+ = A^+$

$$Ax \triangleq b \Rightarrow x \triangleq A^+ b$$

$$A^+ = (A^T A)^{-1} A^T$$

## 2.6 Sensitivity and condition number

matrix with  $\text{rank}(A) = n$

$$\text{cond}(A) = \|A\|_2 \cdot \|A^+\|_2$$

if  $\text{rank}(A) < n \Rightarrow \text{cond}(A) = \infty$

## 2.7 What method to use?

- Easiest  $\Rightarrow$  normal equations, however computationally quite expensive and error proportional to  $(\text{cond}(A))^2$
- Most efficient and accurate orthogonalization method  
 $\Rightarrow$  Householder method.
- SVD most expensive but offers superb robustness

# Eigenvalue Problems

## 3.1 Introduction

$$Ax = \lambda x \Rightarrow \lambda = \text{eigenvalue}$$

$x = \text{eigenvector}$

Characteristic polynomial:

$$\det(A - \lambda I) = 0$$

$\Rightarrow$  Not a good numerical way of computing the eigenvalues.

Properties:

- If  $A$  is symmetric/Hermitian, all are real
- $Ax = \lambda x \Rightarrow (A - \lambda I)x = (A - \lambda I)x$ ,  $x$  remains unchanged
- $A^{-1}$  and  $A$  have same eigenvectors but  $\frac{1}{\lambda}$
- $A^k$  and  $A$  have same eigenvectors but  $\lambda^k$
- polynomials:  $p(t) \Rightarrow p(A) = p(\lambda)x$
- Similarity if there is a  $B = T^{-1}AT$ , it has the same eigenvalues but eigenvectors  $Ty$

## 3.2 Solving eigenvalue problems

Power Iteration

$\Rightarrow$  Estimates the dominant eigenvalue and vector.

Works by multiplying a nonzero vector repeatedly by the matrix.

Usually works well but might fail:

- More than 1 eigenvalue with max. modulus  $\Rightarrow$  converges to a linear combination of corresponding eigenvectors
- Real matrix and starting vector can't converge to a complex vector.

## Inverse iteration.

$\Rightarrow$  Search for the smallest eigenvalue

Using power iteration with LU factorization.

## Rayleigh quotient iteration.

finding  $\lambda$  can be considered a linear least squares problem.

$$x\lambda \approx Ax$$

$$\Rightarrow \lambda = \frac{x^T A x}{x^T x}$$

for an approximate eigenvector  $x$ .

## Deflation.

remove a known eigenvalue from a matrix

so that  $A_1 \Rightarrow U_1^T X = \lambda_1 \Rightarrow A - \lambda_1 U_1^T$  has eigenvalues:

$\Rightarrow$  Not a good way of calculating it.  
 $\lambda = \lambda_2 \dots \lambda_n$

## QR-iteration:

$\Rightarrow$  Fastest and most used method.

$$A_m = Q_m R_m$$

$$A_{m+1} = R_m Q_m$$

$\Rightarrow$  converges to a matrix with the eigenvalues on its diagonal.

## 2.3. Implement

Scipy. linalg. eig.

# 4.1 Numerical Equations

## 4.1.1 Introduction

$$P(x) = g \quad \text{or} \quad P(x) = 0$$

$\Rightarrow g = 0$  the roots.

## 4.1.2 Number of Solutions

possible to have degenerate solutions

$\rightarrow$  multiple roots.

if  $P(x) = P'(x) = P''(x) = \dots = P^{(m-1)}(x) = 0$  and  $P^{(m)}(x) \neq 0$

$x$  has multiplicity  $m$

if  $m=1 \Rightarrow$  simple root.

## 4.1.3 Instability

condition number of  $g$  for root  $x^*$  is  $\|P'(x^*)\|$

for a multiple root is the cond. number in  $g^{-1}$  like

Multiple dimensions: Cond. Number:  $\|\mathcal{J}_g(x)\|^{-1}$

with  $\mathcal{J} = \text{jacobian}$

## 4.4 Convergence rates and skipping

Convergence rate = effectiveness of a certain algorithm

$e_k = x_k - x^*$  error at iteration  $k$

converge with rate  $r$  if:  $\lim_{k \rightarrow \infty} \frac{\|e_{k+1}\|}{\|e_k\|^r} = C$

for  $C > 0$

- $r=1$  and  $C \leq 1$ : linear convergence

- $r > 1$ : superlinear convergence

- $r=2$ : quadratic convergence

- $r=3$ : cubic convergence

Stopping criterion is if

$$\|x_{k+1} - x_k\| / \|x_k\| \leq \epsilon \text{ for error tolerance } \epsilon$$

## Newton's method

### Bisection Method

Begin with a bracket and iteratively reduce its length.  
convergence  $\tau = 1$  and  $c = 95$

Always discards half of the bracket at each iteration  
for a starting interval  $[a, b]$ , after  $k$  iterations is

$$\text{the length equal to } (b-a)/2^k$$

for tolerance  $\epsilon$  it requires  $\log_2(\frac{b-a}{\epsilon})$  iterations

### Fixed-point iteration

$x = g(x)$  is a fixed point

$$x_{k+1} = g(x_k)$$

so its solutions are  $g(x) = 0$ .

if  $|g'(x^*)| < 1$  it is locally convergent

if  $|g'(x^*)| > 1$  then the scheme diverges for every initial value different from  $x^*$

ideally  $|g'(x^*)| = 0 \Rightarrow \lim_{k \rightarrow \infty} \frac{\|x_{k+1}\|}{\|x_k\|} = \frac{|g''(x^*)|}{2}$

### Newton's method

$$f(x+\Delta) \approx f(x) + \Delta f'(x), \Delta = -\frac{f(x)}{f'(x)} \quad \text{if } f'(x) \neq 0$$

$\Rightarrow$  iterative.

- for simple roots it has asymptotic quadratic convergence

- for multiple roots with multiplicity  $m$  (linear conver.)

$$\text{with } C = 1 - \frac{1}{m}$$

## Second Method

$$f'(x_e) = \frac{f(x_e) - f(x_{e-1})}{x_e - x_{e-1}}$$

Use Newton with this as derivative.

## Inverse interpolation

fit a polynomial p to the values of  $x_e$  as function of  $f(x_e)$ , next solution is  $p(0)$ .

Inverse quadratic interpolation  $\Rightarrow$  fits a parabola.

## Rootfinding in Scipy

Brent method  $\Rightarrow$  is a safeguarded method.

## Roots of polynomial functions

find all roots:

- Use one of the methods per root.
- Use a routine.

$\Rightarrow$  Form the companion matrix and compute eigen.

## Solution of nonlinear equations

Define Jacobian:

$J_f(\lambda) \Delta = -f(x)$  then  $x + \Delta$  is an approximate zero of  $f$ .

Able to use Newton and secant.

# 5.1 Optimization

## 5.1.1 Introduction

Some constraints that need to be fulfilled.

concepts and notation

function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  and a set  $S \subseteq \mathbb{R}^n$  we seek  $x^* \in S$  such that  $f$  attains a minimum at  $x^*$

$x^*$  is called a minimizer of the objective function

$S$  is the feasible set, if  $S = \mathbb{R}^n \rightarrow$  unconstrained

Use local minimum instead of global minimum

## 5.2 Optimality conditions

Unconstrained optimality conditions

- gradient:  $\nabla f(x^*) = 0$

- Hessian matrix:

  - Positive definite:  $x^*$  is a minimum

  - Negative definite:  $x^*$  is a maximum

  - Indefinite:  $x^*$  is a saddle point.

## 5.3 Optimization in 1Dimension

A function  $f$  is unimodal on an interval  $[a, b]$  if there is a  $x^*$  that minimizes the function.

**Golden section search**

Compare function values to exclude parts of the interval, the relative distance is the golden ratio, repeat until given error tolerance.

## Successive parabolic interpolation

Use 3 points and fits a parabola,  
minimum of parabola is used as approximate  
of the minimum

16

## Newton's method

$$f(x+h) \approx f(x) + f'(x)h + \frac{1}{2} f''(x)h^2$$

## 5.4 multidimensional optimization

### Direct search

Nelder and mead:  
 $f: \mathbb{R}^m \rightarrow \mathbb{R}$  evaluated  
at  $m+1$  starting points.

A new point replaces the worst point, the point  
is generated on the line of worst point-extremum of the others.  
Useful for nonsmooth objective functions.

### Steepest descent

$$\phi(\alpha) = f(x + \alpha n) \quad \text{with } n = -\nabla f(x)$$

$\Rightarrow$  1 dimensional problem, solvable.

Method is reliable but ~~can have slow progress~~.

### Newton's method

$$f(x+\alpha) \approx f(x) + \nabla f(x)^T \alpha + \frac{1}{2} \alpha^T H f(x) \alpha$$

$$H f(x) \alpha = -\nabla f(x)$$

## Quasi-Newton Methods

$$x_{k+1} = x_k - \alpha_k B_k^{-1} \nabla f(x_k)$$

$\alpha_k$ : line search parameter

$B_k$ : approximation of the Hessian matrix

## Secant updating Methods

a factorization of  $B_k$  updating instead of  $B_k$

## Conjugate gradient method

uitable for large problems.

modifies the gradient at each step to remove components in previous directions.

conjugate product:  $(x_i, y) = x^T H y$

## 5.5 Nonlinear Least Squares

residual:  $r_i(x) = y_i - g_i(x)$

To minimize  $\phi(x) = \frac{1}{2} r(x)^T r(x)$

## Gauss-Newton method

$$(J^T(x_k) J(x_k)) \Delta_k = -J^T(x_k) r(x_k)$$

$$\Rightarrow J(x_k) \Delta_k \approx -r(x_k)$$

## Leverberg-Marquardt Method

When Gauss-Newton yields ill-conditioned or rank-deficient linear least squares,

$$(J^T(x_k) J(x_k) + \mu_k I) \Delta_k = -J^T(x_k) r(x_k)$$

$$\begin{bmatrix} J(x_k) \\ \sqrt{\mu_k} I \end{bmatrix} \Delta_k = \begin{bmatrix} -r(x_k) \\ 0 \end{bmatrix}$$

## 3.6 constrained optimization

18

Search local minima

- Equality constraints:  $g(x) = 0$
- Inequality constraints:  $h(x) \leq 0$

# 4. Interpolation

## 4.1 Interpolation

Interpolation means fitting some function to given data so that the function has the same values as the given data.

for given data  $(t_i, y_i) \Rightarrow f(t_i) = y_i$ : find the interpolant

## 4.2 Polynomial interpolation of discrete data

### Monomial basis

n monomials:  $\phi_j(t) = t^{j-1}$

equation we want to solve:

$$Ax = \begin{bmatrix} 1 & t_1 & t_1^2 & \dots & t_1^{n-1} \\ 1 & t_2 & t_2^2 & \dots & t_2^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & t_n & t_n^2 & \dots & t_n^{n-1} \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} = \begin{bmatrix} y_1 & y_2 & \dots & y_n \end{bmatrix} = y$$

$\hookrightarrow$  Vandermonde matrix

polynomial can be evaluated with Horner's rule or known as nested evaluation or synthetic division.

$$P_{n-1}(t) = x_1 + t(x_2 + t(x_3 + t(\dots(x_{n-1} + t x_n)\dots)))$$

### Lagrange interpolation

for given data  $(t_i, y_i)$ , Lagrange basis function:

$$\ell_j(t) = \frac{\prod_{k=1, k \neq j}^n (t - t_k)}{\prod_{k=1, k \neq j} (t_j - t_k)}$$

polynomial is given by:  $P_{n-1}(t) = y_1 \ell_1(t) + y_2 \ell_2(t) + \dots + y_n \ell_n(t)$

## Newton interpolation

Newton basis functions:

$$\pi_j(t) = \prod_{k=1}^{j-1} (t - t_k)$$

$$P_{n-1}(t) = x_0 + x_1(t-t_1) + x_2(t-t_1)(t-t_2) \dots$$

$$P_{j+1}(t) = p_j(t) + x_{j+1} \pi_{j+1}(t)$$

## Polynomial interpolation of a continuous function

Chebyshev points defined on interval  $[-1, 1]$ :

$$t_i = \cos\left(\frac{(i+1)\pi}{2n}\right)$$

## 6.3. piecewise polynomial interpolation

interval divided into subintervals  
and lower degree polynomials are fitted.

Points at which the interpolations change are knots.

## Cubic spline interpolation

A spline with  $m$  knots  $\Rightarrow$  Cubic Spline has  $4(m-1)$   
free parameters.

Used in software.

## 7.1 Integration

based on Riemann sums:

$$R_n = \sum_{i=1}^n (x_{i+1} - x_i) f(s_i)$$

if  $R_n$  is finite, then the function is Riemann integrable.

## 7.2 Uniqueness, existence, conditioning

- Integrable if it is bounded
- Uniqueness is built into the definition
- Well behaved with a small condition number.

## 7.3 Numerical Quadrature

$$I(f) = \int_a^b f(x) dx = F(b) - F(a)$$

Numerical approximation is numerical quadrature.

quadrature rule:

$$Q_m(f) = \sum_{i=1}^m w_i f(x_i)$$

$x_i$ : nodes,  $w_i$ : weights

to find weights  $\Rightarrow$  method of undetermined coefficients

general system of moment equations:

$$\begin{bmatrix} x_1 & x_2 & \cdots & x_m \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix} = \begin{bmatrix} b - a \\ (b^2 - a^2)/2 \\ \vdots \\ (b^m - a^m)/m \end{bmatrix}$$

accuracy and useful concepts

degree 'd'; if it is exact for polynomial of degree d but not d+1

degree gives accuracy

Rough error bound:

$$\|I(p) - Q_m(f)\| \leq f h^{m+1} \|f''\|_\infty$$

$$h = \max |x_{i+1} - x_i|, i=1 \dots n-1$$

## Newton-Cotes quadrature

open Newton-Cotes rule has nodes:

$$x_i = a + i(b-a)/(n+1)$$

closed:

$$x_i = a + (i+1)(b-a)/(n+1)$$

- Midpoint rule:

$$M(p) = (b-a) \int p\left(\frac{a+b}{2}\right)$$

- Trapezoid rule

$$T(p) = \frac{b-a}{2} (p(a) + p(b))$$

- Simpson's rule

$$S(p) = \frac{b-a}{6} (p(a) + 4p\left(\frac{a+b}{2}\right) + p(b))$$

## Error Analysis

- Midpoint rule:

$$\begin{aligned} I(p) &= p(b-a) + \frac{8}{24} (b-a)^3 + \frac{8}{1920} (b-a)^5 \\ &= M(p) + E(p) + F(p) \end{aligned}$$

- Trapezoid rule:

$$E(p) \approx \frac{T(p) - M(p)}{3}$$

• Simpson's rule:

$$I(p) = \frac{2}{3} M(p) + \frac{1}{3} T(p) - \frac{2}{3} F(p) + \dots = S(p) - \frac{2}{3} F(p)$$

$\Rightarrow$  Newton-Cotes quadrature are relatively easy to derive and apply but has some serious drawbacks

Clebschaw-Curtis quadrature

quadrature based on Chebyshev points.

Gaussian quadrature

nodes and weights are chosen to maximize the degree of the resulting quadrature rule

- nodes are symmetrically placed about the midpoint
- nodes are usually irrational numbers
- more difficult than Newton-Cotes

Composite quadrature

divide interval in subintervals with length  $h = \frac{b-a}{n}$

and taking a simple quadrature in each subinterval

Adaptive quadrature

the interval of integration is selectively refined to reflect the behavior of any particular integrand function.

## 7.4 other integration rules

### Tabular Data

Integrable by piecewise interpolation

### Improper Integrals

- Integrating an integral with unbounded interval of integration:
  - Replace infinities with finites
  - Transform variable so that new interval is finite
  - Use a quadrature rule
- Integrating an integral with singularities
  - Remove by transforming the variable or by dividing out or substracting the singularity.

### Double Integrals

### Multiple Integrals

$\Rightarrow$  Monte carlo

## 7.5 numerical differentiation

### Finite Difference Approximations

- First derivative:

$$\text{- forward difference formula: } f'(x) \approx \frac{f(x+h) - f(x)}{h}$$

$$\text{- backwards difference formula: } f'(x) \approx \frac{f(x) - f(x-h)}{h}$$

$$\text{- central difference formula: } f'(x) \approx \frac{f(x+h) - f(x-h)}{2h}$$

• Second derivative:

- centered:  $f''(x) \approx \frac{f(x+h) - 2f(x) + f(x-h)}{h^2}$

## 7.6 Richardson Extrapolation

Based on values for nonzero step sizes, estimate what the value would be for a step size of zero.

Examples: Romberg integration.

# Unit 2: Numerical Solution of ODE

## 2.1 Introduction

Differential equations has derivatives in it

Only one independent variable : ordinary differential equation.

## 2.1 Numerically Solving ODE

### Euler forward method

Euler's method:  $t_{k+1} = t_k + h_a$ :

$$y_{k+1} = y_k + h_a f(t_k, y_k)$$

$\Rightarrow$  single step method.

- finite difference approximation

$$y'(t) = f(t, y)$$

- Taylor series.

### Accuracy and Stability

- Rounding error

- Truncation error (method error)

- Global error  $\Rightarrow e_k = y_k - \hat{y}(t_k)$

- Local error, error made in one step of the method.

Accuracy is said to be of order p if

$$|e| = O(h_a^{p+1})$$

# Stability

## Implicit methods

Euler backward method:

$$y_{t+1} = y_t + h \cdot f(t_{t+1}, y_{t+1})$$

Stability of euler backward method

to stable:

$$\left| \frac{1}{1-h\lambda} \right| < 1 \quad \text{with } y' = \lambda y$$

## Stiffness

a stable ODE is stiff if its Jacobian has eigenvalues that differ greatly in magnitude.

## Rung-Kutta methods

Taylor series methods take higher derivatives into account

RK-method uses approximations for the higher derivatives

Butcher tables are a way to write down all necessary information in an elegant way.

## Adaptive stepsize

$$h_{\text{opt}} = h_{\text{current}} \left( \frac{\epsilon}{h_{\text{current}} T} \right)^{\frac{1}{N}}$$

classical RK-method uses a fixed step.

Huen's method uses adaptive stepsize as well as Dormand-Prince method

Bogacki-Shampine has the first-name-as-last property

## Extrapolation method

extrapolating a function to fit the computed data.

example is Richardson extrapolation.

## Multi-step methods

uses information at more than one previous point to estimate the solution.

example: Adams-Basforth method.

Properties:

- Use single-step in the beginning
- changing step size is complicated
- local error estimation: predictor - corrector
- Implicit methods have higher stability than explicit
- Implicit can be effective to compute stiff problems.

## Multistep methods

based on polynomial interpolation

## 7.3 BVPs for ODEs

### Introduction

if the side conditions are specified at more than one point  $\Rightarrow$  boundary value problem

## Shooting Method

for 2 point boundary value problem

replaces BVP by a sequence of initial value problems

## 7.1 Introduction

partial differential equation involves partial derivatives of an unknown function with respect to more than one independent variable.

## 7.2 Classification and examples

some important:

- Heat equation
- Wave equation
- Laplace equation

general form:

$$au_{xx} + bu_{xy} + cu_{yy} + du_x + eu_y + fu + g = 0$$

- if:
- $b^2 - 4ac > 0$ : hyperbolic
  - $b^2 - 4ac = 0$ : parabolic
  - $b^2 - 4ac < 0$ : elliptic

## 7.3 Solving time-dependent problems

2 methods:

- Semidiscrete methods: discretize in space but leave the time variable continuous, this then forms an ODE
- Fully discrete methods: all variables are discrete

### Heat equation

$$u_t = c u_{xx} \quad 0 \leq x \leq L, t \geq 0$$

$$\cdot u(0, x) = f(x); 0 \leq x \leq L$$

$$\cdot u(t, 0) = \alpha, u(t, L) = \beta, t \geq 0$$

solve with a semidiscrete method

## Wave equation

$$u_{tt} = c u_{xx}, \quad 0 \leq x \leq L, t \geq 0$$

- $u(0, x) = p(x), \quad u_t(0, x) = q(x), \quad 0 \leq x \leq L$

- $u(t, 0) = \alpha, \quad u(t, L) = \beta, \quad t \geq 0$

Solve with a fully discrete method

### Implicit methods

Some methods have no stability restriction on the relative sizes of  $\Delta t$  and  $\Delta x$ .

Crank-Nicolson Method (Trapezoid method) is unconditionally stable.

## 2.4 Solving the Laplace problem

### Laplace equation

Special case of Poisson equation:

$$u_{xx} + u_{yy} = f(x, y)$$

Laplace B.C.  $f \neq 0$

We have:

- Dirichlet boundary conditions: solution  $u$  is specified.
- Neumann boundary conditions:  $u_x$  or  $u_y$  is specified.
- Robin boundary conditions: a combination of the above.

Laplace equation = potential equation

Solve with a finite difference method

## F1. Motivation

Euler's identity:  $e^{i\theta} = \cos\theta + i\sin\theta$

$$\Rightarrow \cos(2\pi f t) = \frac{e^{2\pi i f t} + e^{-2\pi i f t}}{2}$$

$$\sin(2\pi f t) = \frac{e^{2\pi i f t} - e^{-2\pi i f t}}{2i}$$

$$w_m = \cos\left(\frac{2\pi}{m}\right) - i\sin\left(\frac{2\pi}{m}\right) = e^{-\frac{2\pi i}{m}}$$

## F2. Discrete Fourier Transform

### Definition

a sequence  $x = [x_0 \dots x_{n-1}]^T$ , discrete Fourier transform

$$\Rightarrow y = [y_0 \dots y_{n-1}]^T \Rightarrow y_m = \sum_{k=0}^{n-1} w_m^{mk} x_k$$

$$\Rightarrow y = F_m x$$

$$F_m^{-1} = \frac{1}{m} F_m^H$$

### Frequencies

$$k_e = k_0 + \frac{k}{f_s} \Rightarrow \Delta = (k_e - k_0) f_s$$

$$w_m^{-mk} = \cos\left(\frac{2\pi m k}{m}\right) + i\sin\left(\frac{2\pi m k}{m}\right)$$

$$= \cos\left(\frac{2\pi m}{m} (k_e - k_0) f_s\right) + i\sin\left(\frac{2\pi m}{m} (k_e - k_0) f_s\right)$$

with  $f_s$  the sampling rate

## Negative frequencies

$$w_n^{-m\ell} = w_n^{(m-n)\ell}$$
 arbitrary integer because  $w_n^m = 1$

$$\Rightarrow x_\ell = \frac{1}{m} \sum_{n=0}^{m-1} y_m w_n^{(m-n)\ell}$$

$$\text{or } x_\ell = \frac{1}{m} \sum_{n=0}^{m-1} y_m w_n^{(n-m)\ell}$$

Highest frequency is  $\frac{f_0}{2} \Rightarrow$  Nyquist frequency

## DC component (zero frequency)

lowest frequency = 0

$\Rightarrow$  corresponding coefficient is sum of  $x_\ell$  (the DC component)

## DFT of a real sequence

DFT of a real sequence is in general complex

- $y_0$  is real
- $m \text{ even}: y_{m/2}$  is real
- $m \text{ in odd}: y_{(m+1)/2}$  may be complex

## FFT - ~~algorithm~~

$$D_{m/2} = \text{diag}(1, w_m, \dots, w_m^{(m/2)-1})$$

To apply  $F_m \Rightarrow$  apply  $F_{m/2}$  to its even and odd subsequences

and scale the result by  $\pm D_{m/2}$

$\Rightarrow$  Fast Fourier Transform

## Limitations and extensions

- input frequency needs to be evenly spaced
- input is assumed to be periodic.
- sequence is assumed to be power of 2 in length.

## Performance demo

33

## FFT applications

### Signal processing

- ⇒ to filter out unwanted high-freq. (take the FFT, set high freq to zero and take inverse)
- processing of data which contains multiple periodicities.

### Efficient calculation of convolutions

2 periodic sequences:

$$f \star v = \sum_{k=0}^{m-1} v_k u_{m-k}$$

⇒ equivalent to multiplication by a circulant matrix

⇒ more efficient to use FFT algorithm to transform the inputs to the frequency domain, compute one point wise multiplication and transform the result back to the domain

### Auto correlation

expresses the similarity between a sequence and a delayed copy of itself.

⇒ convolution of the sequence with a reversed copy of itself.

### Fast polynomial multiplication

polynomial product as a convolution:

$$g = f_1 * f_2$$

## 1.1 Pseudo-random number generation

Any computational result is not random.

Don't use LCG for serious applications

### Key ingredients of a PRNG

- A seed
- Algorithm based on internal state
- A recurrence relation.
- Fixed parameters

### Properties of pseudo-random sequences

1. sequence is deterministic
2. a limited number of pseudo-random values
3. pseudo-random sequences must be periodic

### Desirable characteristics of a PRNG

1. uniformly distributed.
2. No apparent statistical correlations between subsequent values

### Modern PRNG algorithms

- Xorshift

### Transformations of univariate continuous distributions.

- Inverse Transform Sampling
- The Box-muller transform

## References

- MersenneTwister
- Xorshift
- Xoshiro256++

Built-in, Numpy and Scipy implementation

Built-in random ⇒ Don't use

Numpy uses permuted Congruential Generator

Scipy has a broad Selection

## 9.2 Basic idea of the Monte Carlo

Relies on this identity:

$$\langle g(X) \rangle = \int g(x) p(x) dx$$

Error estimation

$$\sqrt{\left[ \frac{1}{N} \sum_{i=1}^N g(x_i) \right]} = \frac{1}{\sqrt{N}} \sqrt{g(x_i)}$$

$$\sigma = \sqrt{\frac{\sqrt{g(x_i)}}{N}}$$

$$\sqrt{g(x_i)} \approx \frac{1}{N-1} \sum_{i=1}^N (g(x_i) - \bar{g})^2$$

Application to error propagation

Common Pitfall

for some choices of  $g$ , it becomes infeasible to converge  
the MC estimate

Monte Carlo integration

$$I = \int_{\Omega} g(x) dx$$

$$\rightarrow I = \sqrt{N} \int_{\Omega} g(x) p_{X \sim U(0,1)} X dx$$

## 7.3 Discrete Markov chains

36

### Definition of a discrete Markov chain

A discrete Markov chain is a stochastic process

- ⇒ the probability of observing  $x_{i+1}$  is determined by:
- previous state  $x_i$
  - index  $i$

$$\Rightarrow p_{x_{i+1}}(x_{i+1}) = \int T_{i \rightarrow i+1}(x_{i+1} | x_i) p_{x_i}(x_i) dx$$

- function  $T_{i \rightarrow i+1}$  is not explicitly dependent on  $x_{i-m}$
- if  $T_{i \rightarrow i+1}$  same for all  $i \Rightarrow$  time-homogeneous chain

important property:  $\int T(x_{i+1} | x_i) dx_{i+1} = 1$

### The stationary distribution of a discrete Markov chain

$$\int T(x_2 | x_1) p(x_1) dx = p(x_2)$$

the chain as a random number generator for its corresponding stationary distribution.

### Statistical techniques:

- Stationary distribution does not always exist
- It is not necessarily unique

### Metropolis - Hastings algorithm

- A stationary distribution satisfies the global balance condition: after convolution with the transition probability, the one probability is recovered.

- Detailed balance:

$$T(x_2|x_1)\rho(x_1) = T(x_1|x_2)\rho(x_2)$$

Metropolis-Hastings defines a Markov chain that satisfies detailed balance for any given stationary distribution.

Constructed:

- Generation step

- Acceptance or rejection step

Algorithm:

1. Start with an initial state for which the probability is non-zero

2. Generate a step according to the proposal distribution  $g(x_2|x_1)$

3. Compute the acceptance ratio:

$$AR = \frac{\rho(x_2)}{\rho(x_1)} \frac{g(x_1|x_2)}{g(x_2|x_1)}$$

4. Repeat steps 2 & 3 until sample size is sufficient