

# **Yet Another Introduction to Numerical Methods**

version May 12, 2021

D.V. Fedorov



# Contents

<b>1</b>	<b>Interpolation</b>	<b>1</b>
1.1	Introduction . . . . .	1
1.2	Polynomial interpolation . . . . .	1
1.3	Spline interpolation . . . . .	3
1.3.1	Linear interpolation . . . . .	3
1.3.2	Quadratic spline . . . . .	4
1.3.3	Cubic spline . . . . .	5
1.3.4	Sub-splines . . . . .	8
1.4	Rational function interpolation . . . . .	10
1.5	Multivariate interpolation . . . . .	10
1.5.1	Nearest-neighbor interpolation . . . . .	10
1.5.2	Piecewise-linear interpolation . . . . .	10
1.5.3	Bi-linear interpolation . . . . .	11
<b>2</b>	<b>Systems of linear equations</b>	<b>15</b>
2.1	Introduction . . . . .	15
2.2	Triangular systems . . . . .	16
2.3	Reduction to triangular form . . . . .	16
2.3.1	QR-decomposition . . . . .	17
2.3.2	LU-decomposition . . . . .	23
2.3.3	Cholesky decomposition . . . . .	24
2.4	Determinant of a matrix . . . . .	25
2.5	Matrix inverse . . . . .	25
<b>3</b>	<b>Ordinary least squares problem</b>	<b>27</b>
3.1	Introduction . . . . .	27
3.2	Ordinary least-squares problem . . . . .	27
3.3	Solution via QR-decomposition . . . . .	28
3.4	Ordinary least-squares curve fitting . . . . .	28

3.4.1	Variances and correlations of fitting parameters . . . . .	29
3.5	Singular value decomposition . . . . .	30
<b>4</b>	<b>Eigenvalues and eigenvectors</b>	<b>33</b>
4.1	Introduction . . . . .	33
4.2	Similarity transformations . . . . .	34
4.2.1	Jacobi eigenvalue algorithm . . . . .	34
4.2.2	QR/QL algorithm . . . . .	36
4.3	Eigenvalues of updated matrix . . . . .	36
4.3.1	Rank-1 update . . . . .	37
4.3.2	Symmetric row/column update . . . . .	38
4.3.3	Symmetric rank-2 update . . . . .	39
4.4	Singular Value Decomposition . . . . .	39
<b>5</b>	<b>Ordinary differential equations</b>	<b>43</b>
5.1	Introduction . . . . .	43
5.2	Error estimate . . . . .	44
5.2.1	Runge's principle . . . . .	44
5.2.2	Different orders . . . . .	45
5.3	Runge-Kutta methods . . . . .	46
5.3.1	Embedded methods with error estimates . . . . .	47
5.4	Multistep methods . . . . .	49
5.4.1	Two-step method . . . . .	50
5.4.2	Two-step method with extra evaluation . . . . .	50
5.5	Predictor-corrector methods . . . . .	51
5.5.1	Two-step method with correction . . . . .	51
5.6	Adaptive step-size control . . . . .	52
<b>6</b>	<b>Nonlinear equations</b>	<b>55</b>
6.1	Introduction . . . . .	55
6.2	Newton's method . . . . .	56
6.3	Quasi-Newton methods . . . . .	58
6.3.1	Broyden's method . . . . .	58
<b>7</b>	<b>Numerical integration</b>	<b>61</b>
7.1	Introduction . . . . .	61
7.2	Rectangle and trapezium rules . . . . .	62
7.3	Quadratures with regularly spaced abscissas . . . . .	63
7.3.1	Classical quadratures . . . . .	64
7.4	Quadratures with optimized abscissas . . . . .	65
7.4.1	Gauss quadratures . . . . .	66

7.4.2	Gauss-Kronrod quadratures . . . . .	70
7.5	Adaptive quadratures . . . . .	70
7.6	Variable transformation quadratures . . . . .	71
7.7	Infinite intervals . . . . .	73
<b>8</b>	<b>Minimization and optimization</b>	<b>75</b>
8.1	Introduction . . . . .	75
8.2	Local minimization . . . . .	75
8.2.1	Newton's methods . . . . .	75
8.2.2	Quasi-Newton methods . . . . .	76
8.2.3	Downhill simplex method . . . . .	78
8.3	Global optimization . . . . .	79
8.3.1	Simulated annealing . . . . .	80
8.3.2	Quantum annealing . . . . .	81
8.3.3	Evolutionary algorithms . . . . .	82
8.4	Implementation in C . . . . .	83
<b>9</b>	<b>Monte Carlo integration</b>	<b>87</b>
9.1	Introduction . . . . .	87
9.2	Plain Monte Carlo sampling . . . . .	88
9.3	Importance sampling . . . . .	89
9.4	Stratified sampling . . . . .	90
9.5	Quasi-random (low-discrepancy) sampling . . . . .	91
9.5.1	Van der Corput and Halton sequences . . . . .	92
9.5.2	Lattice rules . . . . .	93
9.6	Implementations . . . . .	93
<b>10</b>	<b>Power iteration methods and Krylov subspaces</b>	<b>97</b>
10.1	Power iteration . . . . .	97
10.2	Inverse iteration . . . . .	98
10.3	Krylov subspaces . . . . .	98
10.4	Arnoldi iteration . . . . .	99
10.5	Lanczos iteration . . . . .	100
10.6	Generalised minimum residual (GMRES) . . . . .	100
<b>11</b>	<b>Fast Fourier transform</b>	<b>103</b>
11.1	Discrete Fourier Transform . . . . .	103
11.1.1	Applications . . . . .	104
11.2	Cooley-Tukey algorithm . . . . .	105
11.3	Multidimensional DFT . . . . .	105



# Chapter 1

## Interpolation

### 1.1 Introduction

In practice one often meets a situation where the function of interest,  $f(x)$ , is only represented by a discrete set of tabulated points,

$$\{x_i, y_i \doteq f(x_i) \mid i = 1 \dots n\},$$

obtained, for example, by sampling, experimentation, or extensive numerical calculations.

*Interpolation* means constructing a (smooth) function, called *interpolating function* or *interpolant*, which passes exactly through the given points and hopefully approximates the tabulated function between the tabulated points. Interpolation is a specific case of *curve fitting* in which the fitting function must go exactly through the data points.

The interpolating function can be used for different practical needs like estimating the tabulated function between the tabulated points and estimating the derivatives/integrals involving the tabulated function.

### 1.2 Polynomial interpolation

Polynomial interpolation uses a polynomial as the interpolating function. Given a table of  $n$  points,  $\{x_i, y_i\}$ , where no two  $x_i$  are the same, one can construct a polynomial  $P(x)$  of the order  $n - 1$  which passes exactly through the points:  $P(x_i) = y_i$ . This

polynomial can be intuitively written in the *Lagrange form*,

$$P(x) = \sum_{i=1}^n y_i \prod_{k \neq i} \frac{x - x_k}{x_i - x_k}. \quad (1.1)$$

The Lagrange interpolating polynomial always exists and is unique.

Table 1.1: Polynomial interpolation in C

```
double polinterp(int n, double *x, double *y, double z) {
  double s=0,p;
  for (int i=0; i<n; i++) {
    p=1; for (int k=0; k<n; k++) if (k!=i) p*=(z-x[k])/(x[i]-x[k]);
    s+=y[i]*p; }
  return s; }
```

Higher order interpolating polynomials are susceptible to the *Runge's phenomenon*: erratic oscillations close to the end-points of the interval, see Figure 1.1.

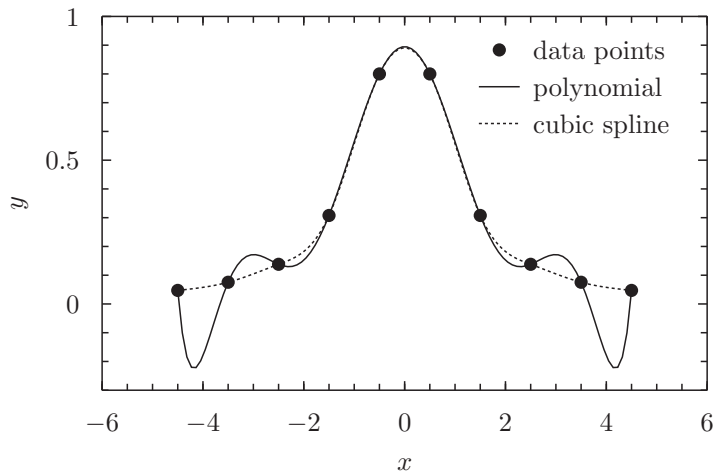


Figure 1.1: Lagrange interpolating polynomial, solid line, showing the Runge's phenomenon: large oscillations at the edges. For comparison the dashed line shows a cubic spline.



## 1.3 Spline interpolation

Spline interpolation uses a *piecewise polynomial*,  $S(x)$ , called *spline*, as the interpolating function,

$$S(x) = s_i(x) \text{ if } x \in [x_i, x_{i+1}] \Big|_{i=1, \dots, n-1}, \quad (1.2)$$

where  $s_i(x)$  is a polynomial of a given order  $k$ . Spline interpolation avoids the problem of Runge's phenomenon. Originally, "spline" was a term for elastic rulers that were bent to pass through a number of predefined points. These were used to make technical drawings by hand.

The spline of the order  $k \geq 1$  can be made continuous at the tabulated points,

$$s_i(x_i) = y_i, \quad s_i(x_{i+1}) = y_{i+1} \Big|_{i=1, \dots, n-1}, \quad (1.3)$$

together with its  $k - 1$  derivatives,

$$\left. \begin{aligned} s'_i(x_{i+1}) &= s'_{i+1}(x_{i+1}), \\ s''_i(x_{i+1}) &= s''_{i+1}(x_{i+1}), \\ &\vdots \\ s_i^{(k-1)}(x_{i+1}) &= s_{i+1}^{(k-1)}(x_{i+1}). \end{aligned} \right|_{i=1, \dots, n-2} \quad (1.4)$$

Continuity conditions (1.3) and (1.4) make  $kn + n - 2k$  linear equations for the  $(n - 1)(k + 1) = kn + n - k - 1$  coefficients of  $n - 1$  polynomials (1.2) of the order  $k$ . The missing  $k - 1$  conditions can be chosen (reasonably) arbitrarily.

The most popular is the cubic spline, where the polynomials  $s_i(x)$  are of third order. The cubic spline is a continuous function together with its first and second derivatives. The cubic spline has a nice feature that it (sort of) minimizes the total curvature of the interpolating function. This makes the cubic splines look good.

Quadratic splines—continuous with the first derivative—are not nearly as good as cubic splines in most respects. In particular they might oscillate unpleasantly when a quick change in the tabulated function is followed by a period where the function is nearly a constant. Cubic splines are somewhat less susceptible to such oscillations.

Linear spline is simply a *polygon* drawn through the tabulated points.

### 1.3.1 Linear interpolation

*Linear interpolation* is a spline with linear polynomials. The continuity conditions (1.3) can be satisfied by choosing the spline in the (intuitive) form

$$s_i(x) = y_i + p_i(x - x_i), \quad (1.5)$$

where

$$p_i = \frac{\Delta y_i}{\Delta x_i}, \quad \Delta y_i \doteq y_{i+1} - y_i, \quad \Delta x_i \doteq x_{i+1} - x_i. \quad (1.6)$$

The linear spline can be easily differentiated,

$$s'_i(x) = p_i, \quad (1.7)$$

and integrated,

$$\int s_i(x) dx = y_i(x - x_i) + p_i \frac{(x - x_i)^2}{2} + \text{const}. \quad (1.8)$$

Linear spline is continuous but its derivative is not.

Table 1.2: Linear interpolation in C

```
#include<assert.h>
double linterp(int n, double* x, double* y, double z){
  assert(n>1 && z>=x[0] && z<=x[n-1]);
  int i=0, j=n-1; /* binary search: */
  while(j-i>1){int m=(i+j)/2; if(z>x[m]) i=m; else j=m;}
  double dy=y[i+1]-y[i], dx=x[i+1]-x[i]; assert(dx>0);
  return y[i]+dy/dx*(z-x[i]);
}
```

Note that the search of the interval  $[x_i \leq x \leq x_{i+1}]$  in an ordered array  $\{x_i\}$  should be done with the *binary search* algorithm (also called *half-interval search*): the point  $x$  is compared to the middle element of the array, if it is less than the middle element, the algorithm repeats its action on the sub-array to the left of the middle element, if it is greater, on the sub-array to the right. When the remaining sub-array is reduced to two elements, the interval is found (see Table 1.2). The average number of operations for a binary search is  $O(\log n)$ .

### 1.3.2 Quadratic spline

Quadratic spline is made of second order polynomials, conveniently written in the form

$$s_i(x) = y_i + p_i(x - x_i) + c_i(x - x_i)(x - x_{i+1}) \Big|_{i=1, \dots, n-1}, \quad (1.9)$$

which identically satisfies the continuity conditions

$$s_i(x_i) = y_i, \quad s_i(x_{i+1}) = y_{i+1} \Big|_{i=1, \dots, n-1}. \quad (1.10)$$

Substituting (1.9) into the derivative continuity condition,

$$s'_i(x_{i+1}) = s'_{i+1}(x_{i+1}) \Big|_{i=1, \dots, n-2}, \quad (1.11)$$

gives  $n - 2$  equations for  $n - 1$  unknown coefficients  $c_i$ ,

$$p_i + c_i \Delta x_i = p_{i+1} - c_{i+1} \Delta x_{i+1} \Big|_{i=1, \dots, n-2}. \quad (1.12)$$

One coefficient can be chosen arbitrarily, for example  $c_1 = 0$ . The other coefficients can now be calculated recursively from (1.12),

$$c_{i+1} = \frac{1}{\Delta x_{i+1}} (p_{i+1} - p_i - c_i \Delta x_i) \Big|_{i=1, \dots, n-2}. \quad (1.13)$$

Alternatively, one can choose  $c_{n-1} = 0$  and make the backward-recursion

$$c_i = \frac{1}{\Delta x_i} (p_{i+1} - p_i - c_{i+1} \Delta x_{i+1}) \Big|_{i=n-2, \dots, 1}. \quad (1.14)$$

In practice, unless you know what your  $c_1$  (or  $c_{n-1}$ ) is, it is better to run both recursions and then average the resulting  $c$ 's. This amounts to first running the forward-recursion from  $c_1 = 0$  and then the backward recursion from  $\frac{1}{2}c_{n-1}$ .

The optimized form (1.9) of the quadratic spline can also be written in the ordinary form suitable for differentiation and integration,

$$s_i(x) = y_i + b_i(x - x_i) + c_i(x - x_i)^2, \text{ where } b_i = p_i - c_i \Delta x_i. \quad (1.15)$$

An implementation of quadratic spline in C is listed in Table 1.3.2

### 1.3.3 Cubic spline

Cubic splines are made of third order polynomials,

$$s_i(x) = y_i + b_i(x - x_i) + c_i(x - x_i)^2 + d_i(x - x_i)^3. \quad (1.16)$$

This form automatically satisfies the first half of continuity conditions (1.3):  $s_i(x_i) = y_i$ . The second half of continuity conditions (1.3),  $s_i(x_{i+1}) = y_{i+1}$ , and the continuity of the first and second derivatives (1.4) give a set of equations,

$$\begin{aligned} y_i + b_i h_i + c_i h_i^2 + d_i h_i^3 &= y_{i+1}, \quad i = 1, \dots, n-1 \\ b_i + 2c_i h_i + 3d_i h_i^2 &= b_{i+1}, \quad i = 1, \dots, n-2 \\ 2c_i + 6d_i h_i &= 2c_{i+1}, \quad i = 1, \dots, n-2 \end{aligned} \quad (1.17)$$

Table 1.3: Quadratic spline in C

```

#include <stdlib.h>
#include <assert.h>
typedef struct {int n; double *x, *y, *b, *c;} qspline;
qspline* qspline_alloc(int n, double* x, double* y){ //builds qspline
qspline *s = (qspline*)malloc(sizeof(qspline)); //spline
s->b = (double*)malloc((n-1)*sizeof(double)); // b_i
s->c = (double*)malloc((n-1)*sizeof(double)); // c_i
s->x = (double*)malloc(n*sizeof(double)); // x_i
s->y = (double*)malloc(n*sizeof(double)); // y_i
s->n = n; for(int i=0;i<n;i++){s->x[i]=x[i]; s->y[i]=y[i];}
int i; double p[n-1], h[n-1]; //VLA from C99
for(i=0;i<n-1;i++){h[i]=x[i+1]-x[i]; p[i]=(y[i+1]-y[i])/h[i];}
s->c[0]=0; //recursion up:
for(i=0;i<n-2;i++)s->c[i+1]=(p[i+1]-p[i]-s->c[i]*h[i])/h[i+1];
s->c[n-2]/=2; //recursion down:
for(i=n-3;i>=0;i--)s->c[i]=(p[i+1]-p[i]-s->c[i+1]*h[i+1])/h[i];
for(i=0;i<n-1;i++)s->b[i]=p[i]-s->c[i]*h[i];
return s; }
double qspline_eval(qspline *s, double z){ //evaluates s(z)
assert(z>=s->x[0] && z<=s->x[s->n-1]);
int i=0, j=s->n-1; //binary search:
while(j-i>1){int m=(i+j)/2; if(z>s->x[m]) i=m; else j=m;}
double h=z-s->x[i];
return s->y[i]+h*(s->b[i]+h*s->c[i]); } //interpolating polynomial
void qspline_free(qspline *s){ //free the allocated memory
free(s->x); free(s->y); free(s->b); free(s->c); free(s); }

```

where  $h_i \doteq x_{i+1} - x_i$ .

The set of equations (1.17) is a set of  $3n-5$  linear equations for the  $3(n-1)$  unknown coefficients  $\{a_i, b_i, c_i \mid i = 1, \dots, n-1\}$ . Therefore two more equations should be added to the set to find the coefficients. If the two extra equations are also linear, the total system is linear and can be easily solved.

The spline is called *natural* if the extra conditions are given as vanishing second derivatives at the end-points,

$$S''(x_1) = S''(x_n) = 0, \quad (1.18)$$

which gives

$$\begin{aligned} c_1 &= 0, \\ c_{n-1} + 3d_{n-1}h_{n-1} &= 0. \end{aligned} \quad (1.19)$$

Solving the first two equations in (1.17) for  $c_i$  and  $d_i$  gives<sup>1</sup>

$$\begin{aligned} c_i h_i &= -2b_i - b_{i+1} + 3p_i, \\ d_i h_i^2 &= b_i + b_{i+1} - 2p_i, \end{aligned} \quad (1.20)$$

where  $p_i \doteq \Delta y_i / h_i$ . The natural conditions (1.19) and the third equation in (1.17) then produce the following tridiagonal system of  $n$  linear equations for the  $n$  coefficients  $b_i$ ,

$$\begin{aligned} 2b_1 + b_2 &= 3p_1, \\ b_i + \left(2\frac{h_i}{h_{i+1}} + 2\right) b_{i+1} + \frac{h_i}{h_{i+1}} b_{i+2} &= 3 \left( p_i + p_{i+1} \frac{h_i}{h_{i+1}} \right) \Big|_{i=1, \dots, n-2}, \\ b_{n-1} + 2b_n &= 3p_{n-1}, \end{aligned} \quad (1.21)$$

or, in the matrix form,

$$\begin{pmatrix} D_1 & Q_1 & 0 & 0 & \cdots \\ 1 & D_2 & Q_2 & 0 & \cdots \\ 0 & 1 & D_3 & Q_3 & \cdots \\ \vdots & \vdots & \ddots & \ddots & \ddots \\ \cdots & \cdots & 0 & 1 & D_n \end{pmatrix} \begin{pmatrix} b_1 \\ \vdots \\ \vdots \\ b_n \end{pmatrix} = \begin{pmatrix} B_1 \\ \vdots \\ \vdots \\ B_n \end{pmatrix} \quad (1.22)$$

where the elements  $D_i$  at the main diagonal are

$$D_1 = 2, \quad D_{i+1} = 2\frac{h_i}{h_{i+1}} + 2 \Big|_{i=1, \dots, n-2}, \quad D_n = 2, \quad (1.23)$$

the elements  $Q_i$  at the above-main diagonal are

$$Q_1 = 1, \quad Q_{i+1} = \frac{h_i}{h_{i+1}} \Big|_{i=1, \dots, n-2}, \quad (1.24)$$

and the right-hand side terms  $B_i$  are

$$B_1 = 3p_1, \quad B_{i+1} = 3 \left( p_i + p_{i+1} \frac{h_i}{h_{i+1}} \right) \Big|_{i=1, \dots, n-2}, \quad B_n = 3p_{n-1}. \quad (1.25)$$

This system can be solved by one run of Gauss elimination and then a run of back-substitution. After a run of Gaussian elimination the system becomes

$$\begin{pmatrix} \tilde{D}_1 & Q_1 & 0 & 0 & \cdots \\ 0 & \tilde{D}_2 & Q_2 & 0 & \cdots \\ 0 & 0 & \tilde{D}_3 & Q_3 & \cdots \\ \vdots & \vdots & \ddots & \ddots & \ddots \\ \cdots & \cdots & 0 & 0 & \tilde{D}_n \end{pmatrix} \begin{pmatrix} b_1 \\ \vdots \\ \vdots \\ b_n \end{pmatrix} = \begin{pmatrix} \tilde{B}_1 \\ \vdots \\ \vdots \\ \tilde{B}_n \end{pmatrix}, \quad (1.26)$$

---

<sup>1</sup>introducing an auxiliary coefficient  $b_n$ .

where

$$\tilde{D}_1 = D_1, \quad \tilde{D}_i = D_i - Q_{i-1}/\tilde{D}_{i-1} \Big|_{i=2,\dots,n}, \quad (1.27)$$

and

$$\tilde{B}_1 = B_1, \quad \tilde{B}_i = B_i - \tilde{B}_{i-1}/\tilde{D}_{i-1} \Big|_{i=2,\dots,n}. \quad (1.28)$$

The triangular system (1.26) can be solved by a run of back-substitution,

$$b_n = \tilde{B}_n/\tilde{D}_n, \quad b_i = (\tilde{B}_i - Q_i b_{i+1})/\tilde{D}_i \Big|_{i=n-1,\dots,1}. \quad (1.29)$$

A C-implementation of cubic spline is listed in Table 1.3.3

### 1.3.4 Sub-splines

Sub-splines are—like splines—piecewise polynomials. However, unlike splines the sub-splines dispense with the demand of maximal differentiability of the spline—hence the name. Instead of maximal differentiability the sub-splines use one (or more) of their free parameters to achieve some other goal. Typically the sub-splines try to minimize the unpleasant wiggles that splines are prone to when the interpolated data contain a sub-set of points that are close to a line.

*Akima sub-spline* [2] is an interpolating function in the form of a piecewise cubic polynomial, similar to the cubic spline,

$$S(x) \Big|_{x \in [x_i, x_{i+1}]} = a_i + b_i(x - x_i) + c_i(x - x_i)^2 + d_i(x - x_i)^3 \doteq A_i(x). \quad (1.30)$$

However, unlike the cubic spline, Akima sub-spline dispenses with the demand of maximal differentiability of the spline—in this case, the continuity of the second derivative—hence the name *sub-spline*. Instead of achieving maximal differentiability Akima sub-splines try to reduce the wiggling which the ordinary splines are typically prone to (see Figure 1.2).

First let us note that the coefficients  $\{a_i, b_i, c_i, d_i\}$  in eq. (1.30) are determined by the values of the derivatives  $S'_i \doteq S'(x_i)$  of the sub-spline through the continuity conditions for the sub-spline and its first derivative,

$$A_i(x_i) = y_i, \quad A'_i(x_i) = S'_i, \quad A_i(x_{i+1}) = y_{i+1}, \quad A'_i(x_{i+1}) = S'_{i+1}. \quad (1.31)$$

Indeed, inserting (1.30) into (1.31) and solving for the coefficients gives

$$a_i = y_i, \quad b_i = S'_i, \quad c_i = \frac{3p_i - 2S'_i - S'_{i+1}}{\Delta x_i}, \quad d_i = \frac{S'_i + S'_{i+1} - 2p_i}{(\Delta x_i)^2}, \quad (1.32)$$

where  $p_i \doteq \Delta y_i / \Delta x_i$ ,  $\Delta y_i \doteq y_{i+1} - y_i$ ,  $\Delta x_i \doteq x_{i+1} - x_i$ .

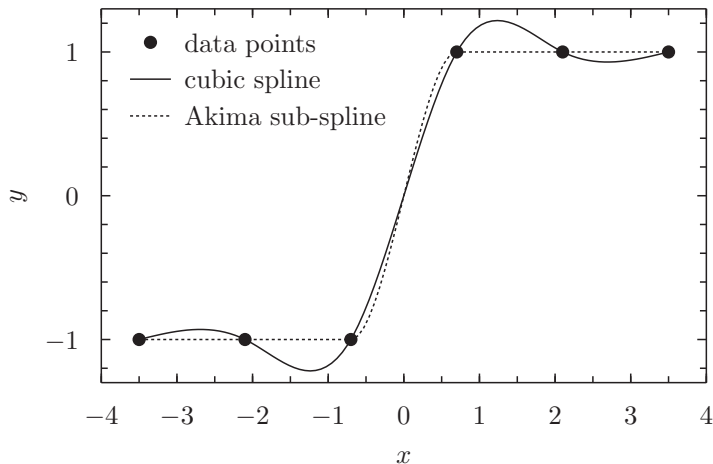


Figure 1.2: A cubic spline (solid line) showing the typical wiggles, compared to the Akima sub-spline (dashed line) where the wiggles are essentially removed.

In the ordinary cubic spline the derivatives  $S'_i$  are determined by the continuity condition of the second derivative of the spline. Sub-splines do without this continuity condition and can instead use the derivatives as free parameters to be chosen to satisfy some other condition.

Akima suggested to minimize the wiggling by choosing the derivatives as linear combinations of the nearest slopes,

$$S'_i = \frac{w_{i+1}p_{i-1} + w_{i-1}p_i}{w_{i+1} + w_{i-1}}, \quad \text{if } w_{i+1} + w_{i-1} \neq 0, \quad (1.33)$$

$$S'_i = \frac{p_{i-1} + p_i}{2}, \quad \text{if } w_{i+1} + w_{i-1} = 0, \quad (1.34)$$

where the weights  $w_i$  are given as

$$w_i = |p_i - p_{i-1}|. \quad (1.35)$$

The idea is that if three points lie close to a line, the sub-spline in this vicinity has to be close to this line. In other words, if  $|p_i - p_{i-1}|$  is small, the nearby derivatives must be close to  $p_i$ .

The first two and the last two points need a special prescription, for example (naively) one can simply use

$$S'_1 = p_1, \quad S'_2 = \frac{1}{2}p_1 + \frac{1}{2}p_2, \quad S'_n = p_{n-1}, \quad S'_{n-1} = \frac{1}{2}p_{n-1} + \frac{1}{2}p_{n-2}. \quad (1.36)$$

Table (1.5) shows a C-implementation of this algorithm.

## 1.4 Rational function interpolation

As the name suggests, the rational interpolation uses a rational function (a ratio of two polynomials) as the interpolant,

$$r_m^k(x) = \frac{p_0 + p_1x + \cdots + p_kx^k}{q_0 + q_1x + \cdots + q_mx^m}, \quad (1.37)$$

where  $m > 0$ . The rational interpolants are (generally) not susceptible to Runge phenomenon and are infinitely differentiable.

Two practical rational interpolants without poles on real axis were suggested by Berrut in 1988 [5],

$$B_1(x) = \frac{\sum_{i=0}^n \frac{(-1)^i}{x - x_i} y_i}{\sum_{i=0}^n \frac{(-1)^i}{x - x_i}}. \quad (1.38)$$

and

$$B_2(x) = \frac{\frac{1}{x - x_0} y_0 + 2 \sum_{i=1}^{n-1} \frac{(-1)^i}{x - x_i} y_i + \frac{(-1)^n}{x - x_n} y_n}{\frac{1}{x - x_0} + 2 \sum_{i=1}^{n-1} \frac{(-1)^i}{x - x_i} + \frac{(-1)^n}{x - x_n}}. \quad (1.39)$$

Berrut interpolants are as slow to evaluate as the Lagrange interpolating polynomial— $O(n)$  operations—but are much more stable.

## 1.5 Multivariate interpolation

Interpolation of a function in more than one variable is called *multivariate interpolation*. The function of interest is represented as a set of discrete points in a multidimensional space. The points may or may not lie on a regular grid.

### 1.5.1 Nearest-neighbor interpolation

Nearest-neighbor interpolation approximates the value of the function at a non-tabulated point by the value at the nearest tabulated point, yielding a piecewise-constant interpolating function. It can be used for both regular and irregular grids.

### 1.5.2 Piecewise-linear interpolation

Piecewise-linear interpolation is used to interpolate functions of two variables tabulated on irregular grids. The tabulated 2D region is triangulated – subdivided into a set of



non-intersecting triangles whose union is the original region. Inside each triangle the interpolating function  $S(x, y)$  is taken in the linear form,

$$S(x, y) = a + bx + cy , \tag{1.40}$$

where the three constants are determined by the three conditions that the interpolating function is equal the tabulated values at the three vertexes of the triangle.

### 1.5.3 Bi-linear interpolation

*Bi-linear interpolation* is used to interpolate functions of two variables tabulated on regular rectilinear 2D grids. The interpolating function  $B(x, y)$  inside each of the grid rectangles is taken as a bilinear function of  $x$  and  $y$ ,

$$B(x, y) = a + bx + cy + dxy , \tag{1.41}$$

where the four constants  $a, b, c, d$  are obtained from the four conditions that the interpolating function is equal the tabulated values at the four nearest tabulated points (which are the vertexes of the given grid rectangle).

Table 1.4: Cubic spline in C

```

#include<stdlib.h>
#include<assert.h>
#include<stdio.h>
typedef struct {int n; double *x,*y,*b,*c,*d;} cubic_spline;
cubic_spline* cubic_spline_alloc(int n, double *x, double *y)
{ // builds natural cubic spline
  cubic_spline* s = (cubic_spline*) malloc(sizeof(cubic_spline));
  s->x = (double*) malloc(n*sizeof(double));
  s->y = (double*) malloc(n*sizeof(double));
  s->b = (double*) malloc(n*sizeof(double));
  s->c = (double*) malloc((n-1)*sizeof(double));
  s->d = (double*) malloc((n-1)*sizeof(double));
  s->n = n; for(int i=0;i<n;i++){s->x[i]=x[i]; s->y[i]=y[i];}
  double h[n-1],p[n-1]; // VLA
  for(int i=0;i<n-1;i++){h[i]=x[i+1]-x[i]; assert(h[i]>0);}
  for(int i=0;i<n-1;i++) p[i]=(y[i+1]-y[i])/h[i];
  double D[n], Q[n-1], B[n]; // building the tridiagonal system:
  D[0]=2; for(int i=0;i<n-2;i++)D[i+1]=2*h[i]/h[i+1]+2; D[n-1]=2;
  Q[0]=1; for(int i=0;i<n-2;i++)Q[i+1]=h[i]/h[i+1];
  for(int i=0;i<n-2;i++)B[i+1]=3*(p[i]+p[i+1]*h[i]/h[i+1]);
  B[0]=3*p[0]; B[n-1]=3*p[n-2]; //Gauss elimination :
  for(int i=1;i<n;i++){ D[i]-=Q[i-1]/D[i-1]; B[i]-=B[i-1]/D[i-1]; }
  s->b[n-1]=B[n-1]/D[n-1]; //back-substitution :
  for(int i=n-2;i>=0;i--) s->b[i]=(B[i]-Q[i]*s->b[i+1])/D[i];
  for(int i=0;i<n-1;i++){
    s->c[i]=(-2*s->b[i]-s->b[i+1]+3*p[i])/h[i];
    s->d[i]=(s->b[i]+s->b[i+1]-2*p[i])/h[i]/h[i];
  }
  return s;
}
double cubic_spline_eval(cubic_spline *s,double z){
  assert(z>=s->x[0] && z<=s->x[s->n-1]);
  int i=0, j=s->n-1; // binary search for the interval for z :
  while(j-i>1){int m=(i+j)/2; if(z>s->x[m]) i=m; else j=m; }
  double h=z-s->x[i]; // calculate the interpolating spline :
  return s->y[i]+h*(s->b[i]+h*(s->c[i]+h*s->d[i]));
}
void cubic_spline_free(cubic_spline *s){ //free the allocated memory
  free(s->x); free(s->y); free(s->b); free(s->c); free(s->d); free(s);}

```

Table 1.5: Akima sub-spline in C

```

#include<assert.h>
#include<stdlib.h>
#include<math.h>
typedef struct {int n; double *x,*y,*b,*c,*d;} akima_spline;
akima_spline* akima_spline_alloc(int n, double *x, double *y){
  assert(n>2); double h[n-1],p[n-1]; /* VLA */
  for(int i=0;i<n-1;i++){h[i]=x[i+1]-x[i]; assert(h[i]>0);}
  for(int i=0;i<n-1;i++) p[i]=(y[i+1]-y[i])/h[i];
  akima_spline *s = (akima_spline*)malloc(sizeof(akima_spline));
  s->x = (double*) malloc(n*sizeof(double));
  s->y = (double*) malloc(n*sizeof(double));
  s->b = (double*) malloc(n*sizeof(double));
  s->c = (double*) malloc((n-1)*sizeof(double));
  s->d = (double*) malloc((n-1)*sizeof(double));
  s->n = n; for(int i=0;i<n;i++){s->x[i]=x[i]; s->y[i]=y[i];}
  s->b[0] =p[0]; s->b[1] =(p[0]+p[1])/2;
  s->b[n-1]=p[n-2]; s->b[n-2]=(p[n-2]+p[n-3])/2;
  for(int i=2;i<n-2;i++){
    double w1=fabs(p[i+1]-p[i]), w2=fabs(p[i-1]-p[i-2]);
    if(w1+w2==0) s->b[i]=(p[i-1]+p[i])/2;
    else s->b[i]=(w1*p[i-1]+w2*p[i])/(w1+w2);
  }
  for(int i=0;i<n-1;i++){
    s->c[i]=(3*p[i]-2*s->b[i]-s->b[i+1])/h[i];
    s->d[i]=(s->b[i+1]+s->b[i]-2*p[i])/h[i]/h[i];
  }
  return s;
}
double akima_spline_eval(akima_spline *s, double z){
  assert(z>=s->x[0] && z<=s->x[s->n-1]);
  int i=0, j=s->n-1;
  while(j-i>1){int m=(i+j)/2; if(z>s->x[m]) i=m; else j=m;}
  double h=z-s->x[i];
  return s->y[i]+h*(s->b[i]+h*(s->c[i]+h*s->d[i]));
}
void akima_spline_free(akima_spline *s){
  free(s->x); free(s->y); free(s->b); free(s->c); free(s->d); free(s);}

```



# Chapter 2

## Systems of linear equations

### 2.1 Introduction

A *system of linear equations* (or *linear system*) is a collection of linear equations involving the same set of unknown variables. A general system of  $n$  linear equations with  $m$  unknowns can be written as

$$\begin{cases} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1m}x_m = b_1 \\ A_{21}x_1 + A_{22}x_2 + \cdots + A_{2m}x_m = b_2 \\ \vdots \qquad \qquad \qquad \vdots \qquad \qquad \qquad \vdots \qquad \qquad \qquad \vdots \\ A_{n1}x_1 + A_{n2}x_2 + \cdots + A_{nm}x_m = b_n \end{cases}, \quad (2.1)$$

where  $x_1, x_2, \dots, x_m$  are the unknown variables,  $A_{11}, A_{12}, \dots, A_{nm}$  are the (constant) coefficients, and  $b_1, b_2, \dots, b_n$  are the (constant) right-hand side terms.

The system can be equivalently written in the matrix form,

$$\mathbf{Ax} = \mathbf{b}, \quad (2.2)$$

where  $\mathbf{A} \doteq \{A_{ij}\}$  is the  $n \times m$  matrix of the coefficients,  $\mathbf{x} \doteq \{x_j\}$  is the size- $m$  column-vector of the unknown variables, and  $\mathbf{b} \doteq \{b_i\}$  is the size- $n$  column-vector of right-hand side terms.

A solution to a linear system is a set of values for the variables  $\mathbf{x}$  which satisfies all equations.

Systems of linear equations occur quite regularly in applied mathematics. Therefore computational algorithms for finding solutions of linear systems are an important part of numerical methods. A system of non-linear equations can often be approximated by a linear system – a helpful technique (called *linearization*) in creating a mathematical model of an otherwise a more complex system.

If  $m = n$  the matrix  $\mathbf{A}$  is called *square*. A square system has a unique solution if  $\mathbf{A}$  is invertible.

## 2.2 Triangular systems

An efficient algorithm to solve numerically a square system of linear equations is to transform the original system into an equivalent *triangular system*,

$$\mathbf{T}\mathbf{y} = \mathbf{c}, \quad (2.3)$$

where  $\mathbf{T}$  is a *triangular matrix* – a special kind of square matrix where the matrix elements either below (upper triangular) or above (lower triangular) the main diagonal are zero.

Indeed, an upper triangular system  $\mathbf{U}\mathbf{y} = \mathbf{c}$  can be easily solved by *back-substitution*,

$$y_i = \frac{1}{U_{ii}} \left( c_i - \sum_{k=i+1}^n U_{ik}y_k \right), \quad i = n, n-1, \dots, 1, \quad (2.4)$$

where one first computes  $y_n = b_n/U_{nn}$ , then substitutes *back* into the previous equation to solve for  $y_{n-1}$ , and repeats through  $y_1$ .

Here is a C-function implementing the in-place<sup>1</sup> back-substitution<sup>2</sup>:

```
void backsub(matrix* U, vector* c){
  for(int i=c->size-1; i>=0; i--){
    double s=vector_get(c, i);
    for(int k=i+1; k<n; k++) s-=matrix_get(U, i, k)*vector_get(c, k);
    vector_set(c, i, s/matrix_get(U, i, i)); } }
```

For a lower triangular system  $\mathbf{L}\mathbf{y} = \mathbf{c}$  the equivalent procedure is called *forward-substitution*,

$$y_i = \frac{1}{L_{ii}} \left( c_i - \sum_{k=1}^{i-1} L_{ik}y_k \right), \quad i = 1, 2, \dots, n. \quad (2.5)$$

## 2.3 Reduction to triangular form

Popular algorithms for reducing a square system of linear equations to a triangular form are *LU-decomposition* and *QR-decomposition*.

<sup>1</sup>here *in-place* means the right-hand side  $\mathbf{c}$  is replaced by the solution  $\mathbf{y}$ .

<sup>2</sup>the functions `vector_get`, `vector_set`, and `matrix_get` are assumed to implement getting and setting the vector- and matrix-elements.

### 2.3.1 QR-decomposition

QR-decomposition is a factorization of a matrix into a product of an orthogonal matrix  $\mathbf{Q}$ , such that  $\mathbf{Q}^T\mathbf{Q} = \mathbf{1}$ , where  $^T$  denotes transposition, and a right triangular matrix  $\mathbf{R}$ , such that

$$\mathbf{A} = \mathbf{QR} . \quad (2.6)$$

QR-decomposition can be used to convert (by multiplying with  $\mathbf{Q}^T$  from the left) a linear system  $\mathbf{Ax} = \mathbf{b}$  into the triangular form,

$$\mathbf{Rx} = \mathbf{Q}^T\mathbf{b} , \quad (2.7)$$

which can be solved directly by back-substitution.

QR-decomposition can also be performed on non-square matrices with few long columns. Generally speaking a rectangular  $n \times m$  matrix  $\mathbf{A}$  can be represented as a product,  $\mathbf{A} = \mathbf{QR}$ , of an orthogonal  $n \times m$  matrix  $\mathbf{Q}$ ,  $\mathbf{Q}^T\mathbf{Q} = \mathbf{1}$ , and a right-triangular  $m \times m$  matrix  $\mathbf{R}$ .

QR-decomposition of a matrix can be computed using several methods, such as Gram-Schmidt orthogonalization, Householder transformation [15], or Givens rotation [11].

#### Gram-Schmidt orthogonalization

*Gram-Schmidt orthogonalization* is an algorithm for orthogonalization of a set of vectors in a given inner product space. It takes a linearly independent set of vectors  $\mathbf{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_m\}$  and generates an orthogonal set  $\mathbf{Q} = \{\mathbf{q}_1, \dots, \mathbf{q}_m\}$  which spans the same subspace as  $\mathbf{A}$ . The algorithm is given as

```

for  $i = 1$  to  $m$  :
   $\mathbf{q}_i \leftarrow \mathbf{a}_i / \|\mathbf{a}_i\|$ 
  for  $j = i + 1$  to  $m$  :  $\mathbf{a}_j \leftarrow \mathbf{a}_j - \langle \mathbf{q}_i | \mathbf{a}_j \rangle \mathbf{q}_i$ 

```

where  $\langle \mathbf{a} | \mathbf{b} \rangle$  is the inner product of two vectors, and  $\|\mathbf{a}\| \doteq \sqrt{\langle \mathbf{a} | \mathbf{a} \rangle}$  is the vector's norm. This variant of the algorithm, where all remaining vectors  $\mathbf{a}_j$  are made orthogonal to  $\mathbf{q}_i$  as soon as the latter is calculated, is considered to be numerically stable and is referred to as *stabilized* or *modified*.

Stabilized Gram-Schmidt orthogonalization can be used to compute QR-decomposition of a matrix  $\mathbf{A}$  by orthogonalization of its column-vectors  $\mathbf{a}_i$  with the inner product

$$\langle \mathbf{a} | \mathbf{b} \rangle = \mathbf{a}^T \mathbf{b} \equiv \sum_{k=1}^n (\mathbf{a})_k (\mathbf{b})_k , \quad (2.8)$$

where  $n$  is the length of column-vectors  $\mathbf{a}$  and  $\mathbf{b}$ , and  $(\mathbf{a})_k$  is the  $k$ th element of the column-vector,

```

for  $i = 1$  to  $m$  :
   $R_{ii} = \sqrt{\mathbf{a}_i^T \mathbf{a}_i}$  ;  $\mathbf{q}_i = \mathbf{a}_i / R_{ii}$ 
  for  $j = i + 1$  to  $m$  :
     $R_{ij} = \mathbf{q}_i^T \mathbf{a}_j$  ;  $\mathbf{a}_j = \mathbf{a}_j - \mathbf{q}_i R_{ij}$  .

```

After orthogonalization the matrices  $\mathbf{Q} = \{\mathbf{q}_1 \dots \mathbf{q}_m\}$  and  $\mathbf{R}$  are the sought orthogonal and right-triangular factors of matrix  $\mathbf{A}$ .

The factorization is unique under requirement that the diagonal elements of  $R$  are positive. For a  $n \times m$  matrix the complexity of the algorithm is  $O(m^2n)$ .

### Gram-Schmidt decomposition with column pivoting

Pivoted decomposition differs from the ordinary Gram-Schmidt in that at each iteration it takes the largest of the remaining columns and thus introduces the permutation matrix  $\mathbf{P}$ ,

$$\mathbf{AP} = \mathbf{QR} , \quad (2.9)$$

that is (generally) chosen so that the diagonal elements of the  $\mathbf{R}$ -matrix are decreasing,

$$|R_{11}| \geq |R_{22}| \geq \dots \geq |R_{mm}| . \quad (2.10)$$

Pivoted QR-decomposition can be used when matrix  $\mathbf{A}$  is rank deficient or its rank is in doubt. With exact arithmetics if  $\text{rank}(\mathbf{A}) = k$  then the sub-matrix of  $\mathbf{R}$  with rows and columns from  $k + 1$  to  $m$  would be zero. Numerical determination of rank requires a criterion for deciding when a small diagonal element of  $\mathbf{R}$  should be treated as zero – a practical choice that depends on both the matrix and the application.

### Householder transformation

A square matrix  $\mathbf{H}$  of the form

$$\mathbf{H} = \mathbf{1} - \frac{2}{\mathbf{u}^T \mathbf{u}} \mathbf{u} \mathbf{u}^T \quad (2.11)$$

is called *Householder matrix*, where the vector  $\mathbf{u}$  is called a *Householder vector*. Householder matrices are symmetric and orthogonal,

$$\mathbf{H}^T = \mathbf{H} , \quad \mathbf{H}^T \mathbf{H} = \mathbf{1} . \quad (2.12)$$

The transformation induced by the Householder matrix on a given vector  $\mathbf{a}$ ,

$$\mathbf{a} \rightarrow \mathbf{H} \mathbf{a} , \quad (2.13)$$

is called a *Householder transformation* or *Householder reflection*. The transformation changes the sign of the affected vector's component in the  $\mathbf{u}$  direction, or, in other words,



makes a reflection of the vector about the hyper-plane perpendicular to  $\mathbf{u}$ , hence the name.

Householder transformation can be used to zero selected components of a given vector  $\mathbf{a}$ . For example, one can zero all components but the first one, such that

$$\mathbf{H}\mathbf{a} = \gamma\mathbf{e}_1, \quad (2.14)$$

where  $\gamma$  is a number and  $\mathbf{e}_1$  is the unit vector in the first direction. The factor  $\gamma$  can be easily calculated,

$$\|\mathbf{a}\|^2 \doteq \mathbf{a}^\top \mathbf{a} = \mathbf{a}^\top \mathbf{H}^\top \mathbf{H} \mathbf{a} = (\gamma\mathbf{e}_1)^\top (\gamma\mathbf{e}_1) = \gamma^2, \quad (2.15)$$

$$\Rightarrow \gamma = \pm \|\mathbf{a}\|. \quad (2.16)$$

To find the Householder vector, we notice that

$$\mathbf{a} = \mathbf{H}^\top \mathbf{H} \mathbf{a} = \mathbf{H}^\top \gamma \mathbf{e}_1 = \gamma \mathbf{e}_1 - \frac{2(\mathbf{u})_1}{\mathbf{u}^\top \mathbf{u}} \mathbf{u}, \quad (2.17)$$

$$\Rightarrow \frac{2(\mathbf{u})_1}{\mathbf{u}^\top \mathbf{u}} \mathbf{u} = \gamma \mathbf{e}_1 - \mathbf{a}, \quad (2.18)$$

where  $(\mathbf{u})_1$  is the first component of the vector  $\mathbf{u}$ . One usually chooses  $(\mathbf{u})_1 = 1$  (for the sake of the possibility to store the other components of the Householder vector in the zeroed elements of the vector  $\mathbf{a}$ ) and stores the factor

$$\frac{2}{\mathbf{u}^\top \mathbf{u}} \equiv \tau \quad (2.19)$$

separately. With this convention one readily finds  $\tau$  from the first component of equation (2.18),

$$\tau = \gamma - (\mathbf{a})_1. \quad (2.20)$$

where  $(\mathbf{a})_1$  is the first element of the vector  $\mathbf{a}$ . For the sake of numerical stability the sign of  $\gamma$  has to be chosen opposite to the sign of  $(\mathbf{a})_1$ ,

$$\gamma = -\text{sign}((\mathbf{a})_1) \|\mathbf{a}\|. \quad (2.21)$$

Finally, the Householder reflection, which zeroes all component of a vector  $\mathbf{a}$  but the first, is given as

$$\mathbf{H} = \mathbf{1} - \tau \mathbf{u} \mathbf{u}^\top, \quad \tau = -\text{sign}((\mathbf{a})_1) \|\mathbf{a}\| - (\mathbf{a})_1, \quad (\mathbf{u})_1 = 1, \quad (\mathbf{u})_{i>1} = -\frac{1}{\tau} (\mathbf{a})_i. \quad (2.22)$$

Now, a QR-decomposition of an  $n \times n$  matrix  $\mathbf{A}$  by Householder transformations can be performed in the following way:

1. Build the size- $n$  Householder vector  $\mathbf{u}_1$  which zeroes the sub-diagonal elements of the first column of matrix  $\mathbf{A}$ , such that

$$\mathbf{H}_1 \mathbf{A} = \left[ \begin{array}{c|ccc} \star & \star & \dots & \star \\ \hline 0 & & & \\ \vdots & & & \\ 0 & & \mathbf{A}_1 & \end{array} \right], \quad (2.23)$$

where  $\mathbf{H}_1 = 1 - \tau_1 \mathbf{u}_1 \mathbf{u}_1^\top$  and where  $\star$  denotes (generally) non-zero matrix elements. In practice one does not build the matrix  $\mathbf{H}_1$  explicitly, but rather calculates the matrix  $\mathbf{H}_1 \mathbf{A}$  in-place, consecutively applying the Householder reflection to columns the matrix  $\mathbf{A}$ , thus avoiding computationally expensive matrix-matrix operations. The zeroed sub-diagonal elements of the first column of the matrix  $\mathbf{A}$  can be used to store the elements of the Householder vector  $\mathbf{u}_1$  while the factor  $\tau_1$  has to be stored separately in a special array. This is the storage scheme used by LAPACK and GSL.

2. Similarly, build the size- $(n - 1)$  Householder vector  $\mathbf{u}_2$  which zeroes the sub-diagonal elements of the first column of matrix  $\mathbf{A}_1$  from eq. (2.23). With the transformation matrix  $\mathbf{H}_2$  defined as

$$\mathbf{H}_2 = \left[ \begin{array}{c|ccc} 1 & 0 & \dots & 0 \\ \hline 0 & & & \\ \vdots & & & \\ 0 & & 1 - \tau_2 \mathbf{u}_2 \mathbf{u}_2^\top & \end{array} \right]. \quad (2.24)$$

the two transformations together zero the sub-diagonal elements of the two first columns of matrix  $\mathbf{A}$ ,

$$\mathbf{H}_2 \mathbf{H}_1 \mathbf{A} = \left[ \begin{array}{cc|ccc} \star & \star & \star & \dots & \star \\ \hline 0 & \star & \star & \dots & \star \\ 0 & 0 & & & \\ \vdots & \vdots & & & \\ 0 & 0 & & \mathbf{A}_3 & \end{array} \right], \quad (2.25)$$

3. Repeating the process zero the sub-diagonal elements of the remaining columns. For column  $k$  the corresponding Householder matrix is

$$\mathbf{H}_k = \left[ \begin{array}{c|c} \mathbf{I}_{k-1} & 0 \\ \hline 0 & 1 - \tau_k \mathbf{u}_k \mathbf{u}_k^\top \end{array} \right], \quad (2.26)$$

where  $\mathbf{I}_{k-1}$  is an identity matrix of size  $k-1$ ,  $\mathbf{u}_k$  is the size- $(n-k+1)$  Householder vector that zeroes the sub-diagonal elements of matrix  $\mathbf{A}_{k-1}$  from the previous step. The corresponding transformation step is

$$\mathbf{H}_k \dots \mathbf{H}_2 \mathbf{H}_1 \mathbf{A} = \left[ \begin{array}{c|c} \mathbf{R}_k & \star \\ \hline 0 & \mathbf{A}_k \end{array} \right], \quad (2.27)$$

where  $\mathbf{R}_k$  is a size- $k$  right-triangular matrix.

After  $n-1$  steps the matrix  $\mathbf{A}$  will be transformed into a right triangular matrix,

$$\mathbf{H}_{n-1} \dots \mathbf{H}_2 \mathbf{H}_1 \mathbf{A} = \mathbf{R}. \quad (2.28)$$

4. Finally, introducing an orthogonal matrix  $\mathbf{Q} = \mathbf{H}_1^\top \mathbf{H}_2^\top \dots \mathbf{H}_{n-1}^\top$  and multiplying eq. (2.28) by  $\mathbf{Q}$  from the left, we get the sought QR-decomposition,

$$\mathbf{A} = \mathbf{Q}\mathbf{R}. \quad (2.29)$$

In practice one does not build explicitly the  $\mathbf{Q}$  matrix but rather applies the successive Householder reflections stored during the decomposition.

### Givens rotations

A Givens rotation is a transformation in the form

$$\mathbf{A} \rightarrow \mathbf{G}(p, q, \theta) \mathbf{A}, \quad (2.30)$$

where  $\mathbf{A}$  is the object to be transformed—matrix of vector—and  $\mathbf{G}(p, q, \theta)$  is the Givens rotation matrix (also known as Jacobi rotation matrix): an orthogonal matrix in the form

$$G(p, q, \theta) = \begin{bmatrix} 1 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \dots & \cos \theta & \dots & \sin \theta & \dots & 0 \\ \vdots & & \vdots & \ddots & \vdots & & \vdots \\ 0 & \dots & -\sin \theta & \dots & \cos \theta & \dots & 0 \\ \vdots & & \vdots & & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 & \dots & 1 \end{bmatrix} \begin{array}{l} \leftarrow \text{row } p \\ \leftarrow \text{row } q \end{array}. \quad (2.31)$$

When a Givens rotation matrix  $\mathbf{G}(p, q, \theta)$  multiplies a vector  $\mathbf{x}$ , only elements  $x_p$  and  $x_q$  are affected. Considering only these two affected elements, the Givens rotation is given explicitly as

$$\begin{bmatrix} x'_p \\ x'_q \end{bmatrix} = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x_p \\ x_q \end{bmatrix} = \begin{bmatrix} x_p \cos \theta + x_q \sin \theta \\ -x_p \sin \theta + x_q \cos \theta \end{bmatrix}. \quad (2.32)$$

Apparently the rotation can zero the element  $x'_q$ , if the angle  $\theta$  is chosen as

$$\tan \theta = \frac{x_q}{x_p} \Rightarrow \theta = \text{atan2}(x_q, x_p). \quad (2.33)$$

A sequence of Givens rotations,

$$\mathbf{G} = \prod_{n \geq q > p=1}^m \mathbf{G}(p, q, \theta_{qp}), \quad (2.34)$$

(where  $n \times m$  is the dimension of the matrix  $\mathbf{A}$ ) can zero all elements of a matrix below the main diagonal if the angles  $\theta_{qp}$  are chosen to zero the elements with indices  $q, p$  of the partially transformed matrix just before applying the matrix  $\mathbf{G}(p, q, \theta_{qp})$ . The resulting matrix is obviously the  $\mathbf{R}$ -matrix of the sought QR-decomposition of the matrix  $\mathbf{A}$  where  $\mathbf{G} = \mathbf{Q}^T$ .

In practice one does not explicitly builds the  $\mathbf{G}$  matrix but rather stores the  $\theta$  angles in the places of the corresponding zeroed elements of the original matrix:

```
#include<gsl/gsl_matrix.h>
#include<math.h>
void givens_qr(gsl_matrix* A){ /* A <- Q,R */
  for(int q=0;q<A->size2;q++){for(int p=q+1;p<A->size1;p++){
    double theta=atan2(gsl_matrix_get(A,p,q),gsl_matrix_get(A,q,q));
    for(int k=q;k<A->size2;k++){
      double xq=gsl_matrix_get(A,q,k), xp=gsl_matrix_get(A,p,k);
      gsl_matrix_set(A,q,k, xq*cos(theta)+xp*sin(theta));
      gsl_matrix_set(A,p,k,-xq*sin(theta)+xp*cos(theta)); }
    gsl_matrix_set(A,p,q,theta); } }
```

When solving the linear system  $\mathbf{Ax} = \mathbf{b}$  one transforms it into the equivalent triangular system  $\mathbf{Rx} = \mathbf{Gb}$  where one calculates  $\mathbf{Gb}$  by successively applying the individual Givens rotations with the stored  $\theta$ -angles:

```
#include<gsl/gsl_vector.h>
#include<gsl/gsl_matrix.h>
#include<math.h>
void givens_qr_QTvec(gsl_matrix* QR, gsl_vector* v){/* v <- Q^T v */
  for(int q=0;q<QR->size2;q++){for(int p=q+1;p<QR->size1;p++){
    double theta = gsl_matrix_get(QR,p,q);
    double vq=gsl_vector_get(v,q), vp=gsl_vector_get(v,p);
    gsl_vector_set(v,q, vq*cos(theta)+vp*sin(theta));
    gsl_vector_set(v,p,-vq*sin(theta)+vp*cos(theta)); } }
```

The triangular system  $\mathbf{Rx} = \mathbf{Gb}$  is then solved by the ordinary back-substitution:

```
#include<gsl/gsl_matrix.h>
#include"givens_qr.h"
void givens_qr_solve(gsl_matrix* QR, gsl_vector* b){
```

```
givens_qr_QTvec(QR, b);
backsub(QR, b); }
```

If one needs to build the  $\mathbf{Q}$ -matrix explicitly, one uses

$$Q_{ij} = \mathbf{e}_i^T \mathbf{Q} \mathbf{e}_j = \mathbf{e}_j^T \mathbf{Q}^T \mathbf{e}_i, \quad (2.35)$$

where  $\mathbf{e}_i$  is the unit vector in the direction  $i$  and where again one can use the successive rotations to calculate  $\mathbf{Q}^T \mathbf{e}_i$ ,

```
#include<gsl/gsl_vector.h>
#include<gsl/gsl_matrix.h>
#include"givens_qr.h"
void givens_qr_unpack_Q(gsl_matrix* QR, gsl_matrix* Q){
  gsl_vector* ei = gsl_vector_alloc(QR->size1);
  for(int i=0; i<QR->size1; i++){
    gsl_vector_set_basis(ei, i);
    givens_qr_QTvec(QR, ei);
    for(int j=0; j<QR->size2; j++){
      gsl_matrix_set(Q, i, j, gsl_vector_get(ei, j)); }
    gsl_vector_free(ei); }
```

Since each Givens rotation only affects two rows of the matrix it is possible to apply a set of rotations in parallel. Givens rotations are also more efficient on sparse matrices.

### 2.3.2 LU-decomposition

LU-decomposition is a factorization of a square matrix  $\mathbf{A}$  into a product of a lower triangular matrix  $\mathbf{L}$  and an upper triangular matrix  $\mathbf{U}$ ,

$$\mathbf{A} = \mathbf{L}\mathbf{U}. \quad (2.36)$$

The linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  after LU-decomposition of the matrix  $\mathbf{A}$  becomes  $\mathbf{L}\mathbf{U}\mathbf{x} = \mathbf{b}$  and can be solved by first solving  $\mathbf{L}\mathbf{y} = \mathbf{b}$  for  $\mathbf{y}$  and then  $\mathbf{U}\mathbf{x} = \mathbf{y}$  for  $\mathbf{x}$  with two runs of forward and backward substitutions.

If  $\mathbf{A}$  is an  $n \times n$  matrix, the condition (2.36) is a set of  $n^2$  equations,

$$\sum_{k=1}^n L_{ik} U_{kj} = A_{ij} \quad |_{i,j=1\dots n}, \quad (2.37)$$

for  $n^2 + n$  unknown elements of the triangular matrices  $\mathbf{L}$  and  $\mathbf{U}$ . The decomposition is thus not unique.

Usually the decomposition is made unique by providing extra  $n$  conditions e.g. by the requirement that the elements of the main diagonal of the matrix  $\mathbf{L}$  are equal one,

$$L_{ii} = 1, \quad i = 1 \dots n. \quad (2.38)$$

The system (2.37) with the extra conditions (2.38) can then be easily solved row after row using the *Doolittle's algorithm*,

```

for  $i = 1 \dots n$  :
   $L_{ii} = 1$ 
  for  $j = i \dots n$  :  $U_{ij} = A_{ij} - \sum_{k < i} L_{ik} U_{kj}$ 
  for  $j = i + 1 \dots n$  :  $L_{ji} = \frac{1}{U_{ii}} \left( A_{ji} - \sum_{k < j} L_{jk} U_{ki} \right)$ 

```

In a slightly different *Crout's algorithm* it is the matrix  $\mathbf{U}$  that has unit diagonal elements,

```

for  $i = 1 \dots n$  :
   $U_{ii} = 1$ 
  for  $j = i \dots n$  :  $L_{ji} = A_{ji} - \sum_{k < i} L_{jk} U_{ki}$ 
  for  $j = i + 1 \dots n$  :  $U_{ij} = \frac{1}{L_{ii}} \left( A_{ji} - \sum_{k < j} L_{jk} U_{ki} \right)$ 

```

Without a proper ordering (permutations) in the matrix, the factorization may fail. For example, it is easy to verify that  $A_{11} = L_{11}U_{11}$ . If  $A_{11} = 0$ , then at least one of  $L_{11}$  and  $U_{11}$  has to be zero, which implies either  $\mathbf{L}$  or  $\mathbf{U}$  is singular, which is impossible if  $\mathbf{A}$  is non-singular. This is however only a procedural problem. It can be removed by simply reordering the rows of  $\mathbf{A}$  so that the first element of the permuted matrix is nonzero (or, even better, the largest in absolute value among all elements of the column below the diagonal). The same problem in subsequent factorization steps can be removed in a similar way. Such algorithm is referred to as *partial pivoting*. It requires an extra integer array to keep track of row permutations.

### 2.3.3 Cholesky decomposition

The Cholesky decomposition of a Hermitian positive-definite matrix  $\mathbf{A}$  is a decomposition in the form

$$\mathbf{A} = \mathbf{L}\mathbf{L}^\dagger, \quad (2.39)$$

where  $\mathbf{L}$  is a lower triangular matrix with real and positive diagonal elements, and  $\mathbf{L}^\dagger$  is the conjugate transpose of  $\mathbf{L}$ .

For real symmetric positive-definite matrices the decomposition reads

$$\mathbf{A} = \mathbf{L}\mathbf{L}^\top, \quad (2.40)$$

where  $\mathbf{L}$  is real.

The decomposition can be calculated using the following in-place algorithm,

$$L_{jj} = \sqrt{A_{jj} - \sum_{k=1}^{j-1} L_{jk}^2}, \quad L_{ij} = \frac{1}{L_{jj}} \left( A_{ij} - \sum_{k=1}^{j-1} L_{ik} L_{jk} \right) \Big|_{i>j}. \quad (2.41)$$

The expression under the square root is always positive if  $\mathbf{A}$  is real and positive-definite.

When applicable, the Cholesky decomposition is about twice as efficient as LU-decomposition for solving systems of linear equations.

## 2.4 Determinant of a matrix

LU- and QR-decompositions allow  $O(n^3)$  calculation of the determinant of a square matrix. Indeed, for the LU-decomposition,

$$\det \mathbf{A} = \det \mathbf{L}\mathbf{U} = \det \mathbf{L} \det \mathbf{U} = \det \mathbf{U} = \prod_{i=1}^n U_{ii} . \quad (2.42)$$

For the Gram-Schmidt QR-decomposition

$$\det \mathbf{A} = \det \mathbf{Q}\mathbf{R} = \det \mathbf{Q} \det \mathbf{R} . \quad (2.43)$$

Since  $\mathbf{Q}$  is an orthogonal matrix  $(\det \mathbf{Q})^2 = 1$ ,

$$|\det \mathbf{A}| = |\det \mathbf{R}| = \left| \prod_{i=1}^n R_{ii} \right| . \quad (2.44)$$

With Gram-Schmidt method one arbitrarily assigns positive sign to diagonal elements of the R-matrix thus removing from the R-matrix the memory of the original sign of the determinant.

However with Givens rotation method the determinant of the individual rotation matrix—and thus the determinant of the total rotation matrix—is equal one, therefore for a square matrix  $\mathbf{A}$  the QR-decomposition  $\mathbf{A} = \mathbf{G}\mathbf{R}$  via Givens rotations allows calculation of the determinant with the correct sign,

$$\det \mathbf{A} = \det \mathbf{R} \equiv \prod_{i=1}^n R_{ii} \quad (2.45)$$

## 2.5 Matrix inverse

The inverse  $\mathbf{A}^{-1}$  of a square  $n \times n$  matrix  $\mathbf{A}$  can be calculated by solving  $n$  linear equations

$$\mathbf{A}\mathbf{x}_i = \mathbf{e}_i \Big|_{i=1, \dots, n} , \quad (2.46)$$

where  $\mathbf{e}_i$  is the unit-vector in the  $i$ -direction: a column where all elements are equal zero except for the element number  $i$  which is equal one. Thus the set of columns

$\{\mathbf{e}_i\}_{i=1,\dots,n}$  form the identity matrix. The matrix made of columns  $\mathbf{x}_i$  is apparently the inverse of  $\mathbf{A}$ .

Here is an implementation of this algorithm using the functions from the Givens rotation chapter,

```
#include<gsl/gsl_vector.h>
#include<gsl/gsl_matrix.h>
#include"givens_qr.h"
void givens_qr_inverse(gsl_matrix* QR, gsl_matrix* B){
    gsl_matrix_set_identity(B);
    for(int i=0; i<QR->size2; i++){
        gsl_vector_view v = gsl_matrix_column(B,i);
        givens_qr_solve(QR,&v.vector); } }
```



## Chapter 3

# Ordinary least squares problem

### 3.1 Introduction

A system of linear equations is considered *overdetermined* if there are more equations than unknown variables. If all equations of an overdetermined system are linearly independent, the system has no exact solution.

An *ordinary least-squares problem* (also called *linear least-squares problem*) is the problem of finding an approximate solution to an overdetermined linear system. It often arises in applications where a theoretical model is fitted to experimental data.

### 3.2 Ordinary least-squares problem

Consider a linear system

$$\mathbf{A}\mathbf{c} = \mathbf{b} , \tag{3.1}$$

where  $\mathbf{A}$  is a  $n \times m$  matrix,  $\mathbf{c}$  is an  $m$ -component vector of unknown variables and  $\mathbf{b}$  is an  $n$ -component vector of the right-hand side terms. If the number of equations  $n$  is larger than the number of unknowns  $m$ , the system is overdetermined and generally has no solution.

However, it is still possible to find an approximate solution — the one where  $\mathbf{A}\mathbf{c}$  is only approximately equal  $\mathbf{b}$  — in the sense that the Euclidean norm of the difference between  $\mathbf{A}\mathbf{c}$  and  $\mathbf{b}$  is minimized,

$$\mathbf{c} : \min_{\mathbf{c}} \|\mathbf{A}\mathbf{c} - \mathbf{b}\|^2 . \tag{3.2}$$

The problem (3.2) is called the ordinary least-squares problem and the vector  $\mathbf{c}$  that minimizes  $\|\mathbf{A}\mathbf{c} - \mathbf{b}\|^2$  is called the *least-squares solution*.

### 3.3 Solution via QR-decomposition

The linear least-squares problem can be solved by QR-decomposition. The matrix  $\mathbf{A}$  is factorized as  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ , where  $\mathbf{Q}$  is  $n \times m$  matrix with orthogonal columns,  $\mathbf{Q}^T\mathbf{Q} = \mathbf{1}$ , and  $\mathbf{R}$  is an  $m \times m$  upper triangular matrix. The Euclidean norm  $\|\mathbf{A}\mathbf{c} - \mathbf{b}\|^2$  can then be rewritten as

$$\begin{aligned} \|\mathbf{A}\mathbf{c} - \mathbf{b}\|^2 &= \|\mathbf{Q}\mathbf{R}\mathbf{c} - \mathbf{b}\|^2 \\ &= \|\mathbf{R}\mathbf{c} - \mathbf{Q}^T\mathbf{b}\|^2 + \|(1 - \mathbf{Q}\mathbf{Q}^T)\mathbf{b}\|^2 \\ &\geq \|(1 - \mathbf{Q}\mathbf{Q}^T)\mathbf{b}\|^2. \end{aligned} \quad (3.3)$$

The term  $\|(1 - \mathbf{Q}\mathbf{Q}^T)\mathbf{b}\|^2$  is independent of the variables  $\mathbf{c}$  and can not be reduced by their variations. However, the term  $\|\mathbf{R}\mathbf{c} - \mathbf{Q}^T\mathbf{b}\|^2$  can be reduced down to zero by solving the  $m \times m$  system of linear equations

$$\mathbf{R}\mathbf{c} = \mathbf{Q}^T\mathbf{b}. \quad (3.4)$$

The system is right-triangular and can be readily solved by back-substitution.

Thus the solution to the ordinary least-squares problem (3.2) is given by the solution of the triangular system (3.4).

### 3.4 Ordinary least-squares curve fitting

Ordinary least-squares curve fitting is a problem of fitting  $n$  (experimental) data points  $\{x_i, y_i \pm \Delta y_i\}_{i=1, \dots, n}$ , where  $\Delta y_i$  are experimental errors, by a linear combination,  $F_{\mathbf{c}}$ , of  $m$  functions  $\{f_k(x)\}_{k=1, \dots, m}$ ,

$$F_{\mathbf{c}}(x) = \sum_{k=1}^m c_k f_k(x), \quad (3.5)$$

where the coefficients  $c_k$  are the fitting parameters.

The objective of the least-squares fit is to minimize the square deviation, called  $\chi^2$ , between the fitting function  $F_{\mathbf{c}}(x)$  and the experimental data [4],

$$\chi^2 = \sum_{i=1}^n \left( \frac{F(x_i) - y_i}{\Delta y_i} \right)^2. \quad (3.6)$$

where the individual deviations from experimental points are weighted with their inverse errors in order to promote contributions from the more precise measurements.

Minimization of  $\chi^2$  with respect to the coefficient  $c_k$  in (3.5) is apparently equivalent to the least-squares problem (3.2) where

$$A_{ik} = \frac{f_k(x_i)}{\Delta y_i}, \quad b_i = \frac{y_i}{\Delta y_i}. \quad (3.7)$$

If  $\mathbf{QR} = \mathbf{A}$  is the QR-decomposition of the matrix  $\mathbf{A}$ , the formal least-squares solution to the fitting problem is

$$\mathbf{c} = \mathbf{R}^{-1} \mathbf{Q}^T \mathbf{b}. \quad (3.8)$$

In practice of course one rather back-substitutes the right-triangular system

$$\mathbf{Rc} = \mathbf{Q}^T \mathbf{b}. \quad (3.9)$$

### 3.4.1 Variances and correlations of fitting parameters

Suppose  $\delta y_i$  is a small deviation of the measured value of the physical observable at hand from its exact value. The corresponding deviation  $\delta c_k$  of the fitting coefficient is then given as

$$\delta c_k = \sum_i \frac{\partial c_k}{\partial y_i} \delta y_i. \quad (3.10)$$

In a good experiment the deviations  $\delta y_i$  are statistically independent and distributed normally with the standard deviations  $\Delta y_i$ . The deviations (3.10) are then also distributed normally with *variances*

$$\langle \delta c_k \delta c_k \rangle = \sum_i \left( \frac{\partial c_k}{\partial y_i} \Delta y_i \right)^2 = \sum_i \left( \frac{\partial c_k}{\partial b_i} \right)^2. \quad (3.11)$$

The standard errors in the fitting coefficients are then given as the square roots of variances,

$$\Delta c_k = \sqrt{\langle \delta c_k \delta c_k \rangle} = \sqrt{\sum_i \left( \frac{\partial c_k}{\partial b_i} \right)^2}. \quad (3.12)$$

The variances are diagonal elements of the *covariance matrix*,  $\Sigma$ , made of *covariances*,

$$\Sigma_{kq} \equiv \langle \delta c_k \delta c_q \rangle = \sum_i \frac{\partial c_k}{\partial b_i} \frac{\partial c_q}{\partial b_i}. \quad (3.13)$$

```

def lsfit (fs:list , x:vector , y:vector , dy:vector) :
    n = x.size ; m = len (fs)
    A = matrix (n,m); b = vector (n)
    for i in range (n):
        b[i] = y[i]/dy[i]
        for k in range (m): A[i,k] = fs [k] (x[i])/dy[i]
    (Q,R) = GramSchmidt (A)
    c = R.backsubstitute (Q.T()*b)
    inverseR = R.backsubstitute (matrix.id_matrix (m))
    S = inverseR*inverseR.T()
    return (c,S)

```

Table 3.1: Python implementation of least-squares curve-fitting by QR-decomposition.

Covariances  $\langle \delta c_k \delta c_q \rangle$  are measures of to what extent the coefficients  $c_k$  and  $c_q$  change together if the measured values  $y_i$  are varied. The normalized covariances,

$$\frac{\langle \delta c_k \delta c_q \rangle}{\sqrt{\langle \delta c_k \delta c_k \rangle \langle \delta c_q \delta c_q \rangle}} \quad (3.14)$$

are called *correlations*.

Using (3.13) and (3.8) the covariance matrix can be calculated as

$$\Sigma = \left( \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \right) \left( \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \right)^{\top} = R^{-1} (R^{-1})^{\top} = (R^{\top} R)^{-1} = (A^{\top} A)^{-1}. \quad (3.15)$$

The square roots of the diagonal elements of this matrix provide the estimates of the errors  $\Delta \mathbf{c}$  of the fitting coefficients,

$$\Delta c_k = \sqrt{\Sigma_{kk}} \Big|_{k=1\dots m}, \quad (3.16)$$

and the (normalized) off-diagonal elements provide the estimates of their correlations.

Table 3.4.1 shows a Python implementation of the ordinary least squares fit via QR decomposition,

An illustration of a fit is shown on Figure 3.1 where a polynomial is fitted to a set of data.

### 3.5 Singular value decomposition

Under the *thin singular value decomposition* we shall understand a representation of a tall  $n \times m$  ( $n > m$ ) matrix  $\mathbf{A}$  in the form

$$\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^{\top}, \quad (3.17)$$

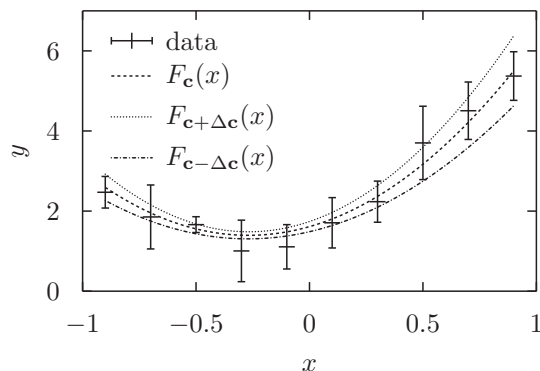


Figure 3.1: Ordinary least squares fit of  $F_{\mathbf{c}}(x) = c_1 + c_2x + c_3x^2$  to a set of data. Shown are fits with optimal coefficients  $\mathbf{c}$  as well as with  $\mathbf{c} + \Delta\mathbf{c}$  and  $\mathbf{c} - \Delta\mathbf{c}$ .

where  $\mathbf{U}$  is an orthogonal  $n \times m$  matrix ( $\mathbf{U}^T\mathbf{U} = \mathbf{1}$ ),  $\mathbf{S}$  is a square  $m \times m$  diagonal matrix with non-negative real numbers on the diagonal (called singular values of matrix  $\mathbf{A}$ ), and  $\mathbf{V}$  is a square  $m \times m$  orthogonal matrix ( $\mathbf{V}^T\mathbf{V} = \mathbf{1}$ ).

Singular value decomposition can be used to solve our linear least squares problem  $\mathbf{A}\mathbf{c} = \mathbf{b}$ . Indeed inserting the decomposition into the equation gives

$$\mathbf{U}\mathbf{S}\mathbf{V}^T\mathbf{c} = \mathbf{b} . \quad (3.18)$$

Multiplying from the left with  $\mathbf{U}^T$  and using the orthogonality of  $\mathbf{U}$  one gets the projected equation

$$\mathbf{S}\mathbf{V}^T\mathbf{c} = \mathbf{U}^T\mathbf{b} . \quad (3.19)$$

This is a square system which can be easily solved first by solving the diagonal system

$$\mathbf{S}\mathbf{y} = \mathbf{U}^T\mathbf{b} \quad (3.20)$$

for  $\mathbf{y}$  and then obtaining  $\mathbf{c}$  as

$$\mathbf{c} = \mathbf{V}\mathbf{y} . \quad (3.21)$$

The covariance matrix (3.15) can be calculated as

$$\mathbf{\Sigma} = (\mathbf{A}^T\mathbf{A})^{-1} = (\mathbf{V}\mathbf{S}^2\mathbf{V}^T)^{-1} = \mathbf{V}\mathbf{S}^{-2}\mathbf{V}^T . \quad (3.22)$$

Singular value decomposition can be found by diagonalising the  $m \times m$  symmetric positive semi-definite matrix  $\mathbf{A}^T\mathbf{A}$  (although this method is not the best for practical calculations, it would do as an educational tool),

$$\mathbf{A}^T\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{V}^T , \quad (3.23)$$

where  $\mathbf{D}$  is a diagonal matrix with eigenvalues of the matrix  $\mathbf{A}^\top \mathbf{A}$  on the diagonal and  $\mathbf{V}$  is the matrix of the corresponding eigenvectors. Indeed it is easy to check that the sought decomposition can be constructed as  $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$  where  $\mathbf{S} = \mathbf{D}^{1/2}$ ,  $\mathbf{U} = \mathbf{A}\mathbf{V}\mathbf{D}^{-1/2}$ .

## Chapter 4

# Eigenvalues and eigenvectors

### 4.1 Introduction

A non-zero column-vector  $\mathbf{v}$  is called the *eigenvector* of a matrix  $\mathbf{A}$  with the *eigenvalue*  $\lambda$ , if

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} . \quad (4.1)$$

If an  $n \times n$  matrix  $\mathbf{A}$  is real and symmetric,  $\mathbf{A}^T = \mathbf{A}$ , then it has  $n$  real eigenvalues  $\lambda_1, \dots, \lambda_n$ , and its (orthogonalized) eigenvectors  $\mathbf{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  form a full basis,

$$\mathbf{V}\mathbf{V}^T = \mathbf{V}^T\mathbf{V} = \mathbf{1} , \quad (4.2)$$

in which the matrix is diagonal,

$$\mathbf{V}^T\mathbf{A}\mathbf{V} = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & & \vdots \\ \vdots & & \ddots & \\ 0 & \cdots & & \lambda_n \end{bmatrix} \equiv \mathbf{D} . \quad (4.3)$$

*Matrix diagonalization* means finding all eigenvalues and (optionally) eigenvectors of a matrix. Once all eigenvalues and eigenvectors are found, the *Eigenvalue Decomposition* of the matrix is given as

$$\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{V}^T . \quad (4.4)$$

Eigenvalues and eigenvectors enjoy a multitude of applications in different branches of science and technology.

## 4.2 Similarity transformations

Orthogonal transformations,

$$\mathbf{A} \rightarrow \mathbf{Q}^T \mathbf{A} \mathbf{Q}, \quad (4.5)$$

where  $\mathbf{Q}^T \mathbf{Q} = \mathbf{1}$ , and, generally, similarity transformations,

$$\mathbf{A} \rightarrow \mathbf{S}^{-1} \mathbf{A} \mathbf{S}, \quad (4.6)$$

preserve eigenvalues and eigenvectors. Therefore one of the strategies to diagonalize a matrix is to apply a sequence of similarity transformations (also called rotations) which (iteratively) turn the matrix into diagonal form.

### 4.2.1 Jacobi eigenvalue algorithm

Jacobi eigenvalue algorithm is an iterative method to calculate the eigenvalues and eigenvectors of a real symmetric matrix by a sequence of Jacobi rotations.

Jacobi rotation is an orthogonal transformation which zeroes a pair of the off-diagonal elements of a (real symmetric) matrix  $\mathbf{A}$ ,

$$\mathbf{A} \rightarrow \mathbf{A}' = \mathbf{J}(p, q)^T \mathbf{A} \mathbf{J}(p, q) : A'_{pq} = A'_{qp} = 0. \quad (4.7)$$

The orthogonal matrix  $\mathbf{J}(p, q)$  which eliminates the element  $A_{pq}$  is called the Jacobi rotation matrix. It is equal identity matrix except for the four elements with indices  $pp$ ,  $pq$ ,  $qp$ , and  $qq$ ,

$$\mathbf{J}(p, q) = \begin{bmatrix} 1 & & & & & & & & & & \\ & \ddots & & & & & & & & & \\ & & \cos \phi & \cdots & \sin \phi & & & & & & 0 \\ & & \vdots & \ddots & \vdots & & & & & & \\ & & -\sin \phi & \cdots & \cos \phi & & & & & & \\ & & & & & & & & \ddots & & \\ 0 & & & & & & & & & & \\ & & & & & & & & & & 1 \end{bmatrix} \begin{matrix} \\ \\ \leftarrow \text{row } p \\ \\ \leftarrow \text{row } q \\ \\ \end{matrix} \quad (4.8)$$

Or explicitly,

$$\begin{aligned} J(p, q)_{ij} &= \delta_{ij} \quad \forall ij \notin \{pq, qp, pp, qq\}; \\ J(p, q)_{pp} &= \cos \phi = J(p, q)_{qq}; \\ J(p, q)_{pq} &= \sin \phi = -J(p, q)_{qp}. \end{aligned} \quad (4.9)$$



After a Jacobi rotation,  $\mathbf{A} \rightarrow \mathbf{A}' = \mathbf{J}^\top \mathbf{A} \mathbf{J}$ , the matrix elements of  $\mathbf{A}'$  become

$$\begin{aligned}
A'_{ij} &= A_{ij} \quad \forall i \neq p, q \wedge j \neq p, q \\
A'_{pi} &= A'_{ip} = cA_{pi} - sA_{qi} \quad \forall i \neq p, q ; \\
A'_{qi} &= A'_{iq} = sA_{pi} + cA_{qi} \quad \forall i \neq p, q ; \\
A'_{pp} &= c^2 A_{pp} - 2scA_{pq} + s^2 A_{qq} ; \\
A'_{qq} &= s^2 A_{pp} + 2scA_{pq} + c^2 A_{qq} ; \\
A'_{pq} &= A'_{qp} = sc(A_{pp} - A_{qq}) + (c^2 - s^2)A_{pq} , \tag{4.10}
\end{aligned}$$

where  $c \equiv \cos \phi$ ,  $s \equiv \sin \phi$ . The angle  $\phi$  is chosen such that after rotation the matrix element  $A'_{pq}$  is zeroed,

$$\tan(2\phi) = \frac{2A_{pq}}{A_{qq} - A_{pp}} \Rightarrow A'_{pq} = 0 , , \phi = \text{atan2}(2A_{pq}, A_{qq} - A_{pp}) . \tag{4.11}$$

A side effect of zeroing a given off-diagonal element  $A_{pq}$  by a Jacobi rotation is that other off-diagonal elements are changed. Namely, the elements of the rows and columns with indices  $p$  and  $q$ . However, after the Jacobi rotation the sum of squares of all off-diagonal elements is reduced. The algorithm repeatedly performs rotations until the off-diagonal elements become sufficiently small.

The convergence of the Jacobi method can be proved for two strategies for choosing the order in which the elements are zeroed:

1. *Classical method*: with each rotation the largest of the remaining off-diagonal elements is zeroed.
2. *Cyclic method*: the off-diagonal elements are zeroed in strict order, e.g. row after row.

Although the classical method allows the least number of rotations, it is typically slower than the cyclic method since searching for the largest element is an  $O(n^2)$  operation. The count can be reduced by keeping an additional array with indexes of the largest elements in each row. Updating this array after each rotation is only an  $O(n)$  operation.

A *sweep* is a sequence of Jacobi rotations applied to all non-diagonal elements. Typically the method converges after a small number of sweeps. The operation count is  $O(n)$  for a Jacobi rotation and  $O(n^3)$  for a sweep.

The typical convergence criterion is that the diagonal elements have not changed after a sweep. Other criteria can also be used, like the sum of absolute values of the off-diagonal elements is small,  $\sum_{i < j} |A_{ij}| < \epsilon$ , where  $\epsilon$  is the required accuracy, or the largest off-diagonal element is small,  $\max |A_{i < j}| < \epsilon$ .

The eigenvectors can be calculated as  $V = \mathbf{1}J_1J_2\dots$ , where  $J_i$  are the successive Jacobi matrices. At each stage the transformation is

$$\begin{aligned} V_{ij} &\rightarrow V_{ij}, j \neq p, q \\ V_{ip} &\rightarrow cV_{ip} - sV_{iq} \\ V_{iq} &\rightarrow sV_{ip} + cV_{iq} \end{aligned} \tag{4.12}$$

Alternatively, if only one (or few) eigenvector  $\mathbf{v}_k$  is needed, one can instead solve the (singular) system  $(A - \lambda_k)\mathbf{v} = 0$ .

### 4.2.2 QR/QL algorithm

An orthogonal transformation of a real symmetric matrix,  $\mathbf{A} \rightarrow \mathbf{Q}^T \mathbf{A} \mathbf{Q} = \mathbf{R} \mathbf{Q}$ , where  $\mathbf{Q}$  is from the QR-decomposition of  $\mathbf{A}$ , partly turns the matrix  $\mathbf{A}$  into diagonal form. Successive iterations eventually make it diagonal. If there are degenerate eigenvalues there will be a corresponding block-diagonal sub-matrix.

For convergence properties it is of advantage to use *shifts*: instead of  $\text{QR}[\mathbf{A}]$  we do  $\text{QR}[\mathbf{A} - s\mathbf{1}]$  and then  $\mathbf{A} \rightarrow \mathbf{R} \mathbf{Q} + s\mathbf{1}$ . The shift  $s$  can be chosen as  $\mathbf{A}_{nn}$ . As soon as an eigenvalue is found the matrix is deflated, that is, the corresponding row and column are crossed out.

Accumulating the successive transformation matrices  $\mathbf{Q}_i$  into the total matrix  $\mathbf{Q} = \mathbf{Q}_1 \dots \mathbf{Q}_N$ , such that  $\mathbf{Q}^T \mathbf{A} \mathbf{Q} = \Lambda$ , gives the eigenvectors as columns of the  $\mathbf{Q}$  matrix.

If only one (or few) eigenvector  $\mathbf{v}_k$  is needed one can instead solve the (singular) system  $(\mathbf{A} - \lambda_k)\mathbf{v} = 0$ .

#### Tridiagonalization.

Each iteration of the QR/QL algorithm is an  $O(n^3)$  operation. On a tridiagonal matrix it is only  $O(n)$ . Therefore the effective strategy is first to make the matrix tridiagonal and then apply the QR/QL algorithm. Tridiagonalization of a matrix is a non-iterative operation with a fixed number of steps.

## 4.3 Eigenvalues of updated matrix

In practice it happens quite often that the matrix  $A$  to be diagonalized is given in the form of a diagonal matrix,  $\mathbf{D}$ , plus an update matrix,  $\mathbf{W}$ ,

$$\mathbf{A} = \mathbf{D} + \mathbf{W}, \tag{4.13}$$

where the update  $\mathbf{W}$  is a simpler, in a certain sense, matrix which allows a more efficient calculation of the updated eigenvalues, as compared to general diagonalization algorithms.

The most common updates are

- symmetric rank-1 update,

$$\mathbf{W} = \mathbf{u}\mathbf{u}^T, \quad (4.14)$$

where  $\mathbf{u}$  is a column-vector;

- symmetric rank-2 update,

$$\mathbf{W} = \mathbf{u}\mathbf{v}^T + \mathbf{v}\mathbf{u}^T; \quad (4.15)$$

- symmetric row/column update – a special case of rank-2 update,

$$\mathbf{W} = \begin{bmatrix} 0 & \dots & u_1 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ u_1 & \dots & u_p & \dots & u_n \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & u_n & \dots & 0 \end{bmatrix} \equiv \mathbf{e}(p)\mathbf{u}^T + \mathbf{u}\mathbf{e}(p)^T, \quad (4.16)$$

where  $\mathbf{e}(p)$  is the unit vector in the  $p$ -direction.

### 4.3.1 Rank-1 update

We assume that a size- $n$  real symmetric matrix  $A$  to be diagonalized is given in the form of a diagonal matrix plus a rank-1 update,

$$\mathbf{A} = \mathbf{D} + \sigma\mathbf{u}\mathbf{u}^T, \quad (4.17)$$

where  $\mathbf{D}$  is a diagonal matrix with diagonal elements  $\{d_1, \dots, d_n\}$  and  $\mathbf{u}$  is a given vector. The diagonalization of such matrix can be done in  $O(m^2)$  operations, where  $m \leq n$  is the number of non-zero elements in the update vector  $\mathbf{u}$ , as compared to  $O(n^3)$  operations for a general diagonalization [12].

The eigenvalue equation for the updated matrix reads

$$(\mathbf{D} + \sigma\mathbf{u}\mathbf{u}^T)\mathbf{q} = \lambda\mathbf{q}, \quad (4.18)$$

where  $\lambda$  is an eigenvalue and  $\mathbf{q}$  is the corresponding eigenvector. The equation can be rewritten as

$$(\mathbf{D} - \lambda\mathbf{1})\mathbf{q} + \sigma\mathbf{u}\mathbf{u}^T\mathbf{q} = 0. \quad (4.19)$$

Multiplying from the left with  $\mathbf{u}^T(\mathbf{D} - \lambda\mathbf{1})^{-1}$  gives

$$\mathbf{u}^T\mathbf{q} + \mathbf{u}^T(\mathbf{D} - \lambda\mathbf{1})^{-1}\sigma\mathbf{u}\mathbf{u}^T\mathbf{q} = 0. \quad (4.20)$$

Finally, dividing by  $\mathbf{u}^T \mathbf{q}$  leads to the (scalar) *secular equation* (or *characteristic equation*) in  $\lambda$ ,

$$1 + \sum_{i=1}^m \frac{\sigma u_i^2}{d_i - \lambda} = 0, \quad (4.21)$$

where the summation index counts the  $m$  non-zero components of the update vector  $\mathbf{u}$ . The  $m$  roots of this equation determine the (updated) eigenvalues<sup>1</sup>.

Finding a root of a rational function requires an iterative technique, such as the Newton-Raphson method. Therefore diagonalization of an updated matrix is still an iterative procedure. However, each root can be found in  $O(1)$  iterations, each iteration requiring  $O(m)$  operations. Therefore the iterative part of this algorithm — finding all  $m$  roots — needs  $O(m^2)$  operations.

Finding roots of this particular secular equation can be simplified by utilizing the fact that its roots are bounded by the eigenvalues  $d_i$  of the matrix  $\mathbf{D}$ . Indeed if we denote the roots as  $\lambda_1, \lambda_2, \dots, \lambda_n$  and assume that  $\lambda_i \leq \lambda_{i+1}$  and  $d_i \leq d_{i+1}$ , it can be shown that

1. if  $\sigma \geq 0$ ,

$$d_i \leq \lambda_i \leq d_{i+1}, \quad i = 1, \dots, n-1, \quad (4.22)$$

$$d_n \leq \lambda_n \leq d_n + \sigma \mathbf{u}^T \mathbf{u}; \quad (4.23)$$

2. if  $\sigma \leq 0$ ,

$$d_{i-1} \leq \lambda_i \leq d_i, \quad i = 2, \dots, n, \quad (4.24)$$

$$d_1 + \sigma \mathbf{u}^T \mathbf{u} \leq \lambda_1 \leq d_1. \quad (4.25)$$

### 4.3.2 Symmetric row/column update

The matrix  $\mathbf{A}$  to be diagonalized is given in the form

$$\mathbf{A} = \mathbf{D} + \mathbf{e}(p)\mathbf{u}^T + \mathbf{u}\mathbf{e}(p)^T = \begin{bmatrix} d_1 & \dots & u_1 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ u_1 & \dots & d_p & \dots & u_n \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & u_n & \dots & d_n \end{bmatrix}, \quad (4.26)$$

where  $\mathbf{D}$  is a diagonal matrix with diagonal elements  $\{d_i | i = 1, \dots, n\}$ ,  $\mathbf{e}(p)$  is the unit vector in the  $p$ -direction, and  $\mathbf{u}$  is a given update vector where the  $p$ -th element can be

<sup>1</sup>Multiplying this equation by  $\prod_{i=1}^m (d_i - \lambda)$  leads to an equivalent polynomial equation of the order  $m$ , which has exactly  $m$  roots.

assumed to equal zero,  $u_p = 0$ , without loss of generality. Indeed, if the element is not zero, one can simply redefine  $d_p \rightarrow d_p + 2u_p$ ,  $u_p \rightarrow 0$ .

The eigenvalue equation for matrix  $\mathbf{A}$  is given as

$$(D - \lambda)\mathbf{x} + \mathbf{e}(p)\mathbf{u}^T\mathbf{x} + \mathbf{u}\mathbf{e}(p)^T\mathbf{x} = 0, \quad (4.27)$$

where  $\mathbf{x}$  is an eigenvector and  $\lambda$  is the corresponding eigenvalue. The component number  $p$  of this vector-equation reads

$$(d_p - \lambda)x_p + \mathbf{u}^T\mathbf{x} = 0, \quad (4.28)$$

while the component number  $k \neq p$  reads

$$(d_k - \lambda)x_k + u_k x_p = 0, \quad (4.29)$$

Dividing the last equation by  $(d_k - \lambda)$ , multiplying from the left with  $\sum_{k=1}^n u_k$ , substituting  $\mathbf{u}^T\mathbf{x}$  using equation (4.28) and dividing by  $x_p$  gives the secular equation,

$$-(d_p - \lambda) + \sum_{k \neq p}^n \frac{u_k^2}{d_k - \lambda} = 0, \quad (4.30)$$

which determines the updated eigenvalues.

### 4.3.3 Symmetric rank-2 update

A symmetric rank-2 update can be represented as two consecutive rank-1 updates,

$$\mathbf{u}\mathbf{v}^T + \mathbf{v}\mathbf{u}^T = \mathbf{a}\mathbf{a}^T - \mathbf{b}\mathbf{b}^T, \quad (4.31)$$

where

$$\mathbf{a} = \frac{1}{\sqrt{2}}(\mathbf{u} + \mathbf{v}), \quad \mathbf{b} = \frac{1}{\sqrt{2}}(\mathbf{u} - \mathbf{v}). \quad (4.32)$$

The eigenvalues can then be found by applying the rank-1 update method twice.

## 4.4 Singular Value Decomposition

Singular Value Decomposition (SVD) is a factorization of matrix  $\mathbf{A}$  in the form

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T, \quad (4.33)$$

where  $\mathbf{D}$  is a diagonal matrix, and  $\mathbf{U}$  and  $\mathbf{V}$  are orthogonal matrices ( $\mathbf{U}^T\mathbf{U} = \mathbf{1}$  and  $\mathbf{V}^T\mathbf{V} = \mathbf{1}$ ).

The elements of the diagonal matrix  $\mathbf{D}$  are called the *singular values* of matrix  $\mathbf{A}$ . Singular values can always be chosen non-negative by changing the signs of the corresponding columns of matrix  $\mathbf{U}$ .

Singular values are equal the square roots of the eigenvalues of the real symmetric matrix  $\mathbf{A}^\top \mathbf{A}$ . Indeed, by construction the matrix  $\mathbf{A}^\top \mathbf{A}$  is positive definite, therefore its eigenvalues are not-negative and its eigenvalue decomposition can be written as

$$\mathbf{A}^\top \mathbf{A} = \mathbf{V} \mathbf{S}^2 \mathbf{V}^\top \quad (4.34)$$

where  $\mathbf{S}$  is the diagonal matrix of square roots of  $\mathbf{A}^\top \mathbf{A}$ , and  $\mathbf{V}$  is the matrix of corresponding eigenvectors. Then the singular value decomposition of matrix  $\mathbf{A}$  can be formally written as

$$\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^\top, \quad \mathbf{U} = \mathbf{A} \mathbf{V} \mathbf{S}^{-1}. \quad (4.35)$$

One algorithm to perform SVD is the *two-sided Jacobi SVD algorithm* which is a generalization of the Jacobi eigenvalue algorithm. In the two-sided Jacobi SVD algorithm one first applies a Givens rotation to symmetrize a pair of off-diagonal elements of the matrix and then applies a Jacobi transformation to eliminate these off-diagonal elements.

It is an iterative procedure where one starts with  $\mathbf{A}_0 = \mathbf{A}$  and then iterates

$$\mathbf{A}_k \rightarrow \mathbf{A}_{k+1} = \mathbf{J}_k^\top \mathbf{G}_k^\top \mathbf{A}_k \mathbf{J}_k. \quad (4.36)$$

Just like in the Jacobi eigenvalue algorithm the iterations are performed in cyclic sweeps over all non-diagonal elements of the matrix. At each iteration the matrix  $\mathbf{G}$  equalizes the corresponding non-diagonal elements, and then the Jacobi transformation zeroes them. The iteration procedure stops when the diagonal elements remain unchanged for a whole sweep.

For a  $2 \times 2$  matrix the two-sided Jacobi SVD transformation is given as following: first, one applies a Givens rotation to symmetrize two off-diagonal elements,

$$\mathbf{A}_k \equiv \begin{bmatrix} w & x \\ y & z \end{bmatrix} \rightarrow \mathbf{A}'_k = \mathbf{G}_k^\top \mathbf{A}_k = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} \begin{bmatrix} w & x \\ y & z \end{bmatrix} = \begin{bmatrix} a & b \\ b & c \end{bmatrix}, \quad (4.37)$$

where the rotation angle  $\theta = \text{atan2}(x - y, w + z)$ ; and, second, one makes the usual Jacobi transformation to eliminate the off-diagonal elements,

$$\begin{aligned} \mathbf{A}'_k &\rightarrow \mathbf{A}_{k+1} = \mathbf{J}_k^\top \mathbf{A}'_k \mathbf{J}_k \\ &= \begin{bmatrix} \cos(\phi) & -\sin(\phi) \\ \sin(\phi) & \cos(\phi) \end{bmatrix} \begin{bmatrix} a & b \\ b & c \end{bmatrix} \begin{bmatrix} \cos(\phi) & \sin(\phi) \\ -\sin(\phi) & \cos(\phi) \end{bmatrix} = \begin{bmatrix} d_1 & 0 \\ 0 & d_1 \end{bmatrix}. \end{aligned} \quad (4.38)$$

The matrices  $U$  and  $V$  are accumulated (from identity matrices) as

$$\mathbf{U}_{k+1} = \mathbf{U}_k \mathbf{G}_k \mathbf{J}_k, \quad (4.39)$$

$$\mathbf{V}_{k+1} = \mathbf{V}_k \mathbf{J}_k. \quad (4.40)$$

If the matrix  $\mathbf{A}$  is a tall  $n \times m$  non-square matrix ( $n > m$ ), the first step should be the QR-decomposition,

$$\mathbf{A} = \mathbf{QR}, \quad (4.41)$$

where  $\mathbf{Q}$  is the  $n \times m$  orthogonal matrix and  $\mathbf{R}$  is a square triangular  $m \times m$  matrix.

The second step is the normal SVD of the square matrix  $\mathbf{R}$ ,

$$\mathbf{R} = \mathbf{U}' \mathbf{D} \mathbf{V}'^T. \quad (4.42)$$

Now the SVD of the original matrix  $\mathbf{A}$  is given as

$$\mathbf{A} = \mathbf{UDV}'^T, \quad (4.43)$$

where

$$\mathbf{U} = \mathbf{QU}' . \quad (4.44)$$

Table 4.1: Jacobi diagonalization in C using `gsl_matrix` and `gsl_vector` as containers.

```

#include<math.h>
#include<gsl/gsl_matrix.h>
#include<gsl/gsl_vector.h>
int jacobi(gsl_matrix* A, gsl_vector* e, gsl_matrix* V){
/* Jacobi diagonalization; upper triangle of A is destroyed;
   e and V accumulate eigenvalues and eigenvectors */
int changed, sweeps=0, n=A->size1;
for (int i=0;i<n;i++)gsl_vector_set(e,i,gsl_matrix_get(A,i,i));
gsl_matrix_set_identity(V);
do{ changed=0; sweeps++; int p,q;
  for (p=0;p<n;p++)for (q=p+1;q<n;q++){
    double app=gsl_vector_get(e,p);
    double aqq=gsl_vector_get(e,q);
    double apq=gsl_matrix_get(A,p,q);
    double phi=0.5*atan2(2*apq,aqq-app);
    double c = cos(phi), s = sin(phi);
    double app1=c*c*app-2*s*c*apq+s*s*aqq;
    double aqq1=s*s*app+2*s*c*apq+c*c*aqq;
    if (app1!=app || aqq1!=aqq){ changed=1;
      gsl_vector_set(e,p,app1);
      gsl_vector_set(e,q,aqq1);
      gsl_matrix_set(A,p,q,0.0);
      for (int i=0;i<p;i++){
        double aip=gsl_matrix_get(A,i,p);
        double aiq=gsl_matrix_get(A,i,q);
        gsl_matrix_set(A,i,p,c*aip-s*aiq);
        gsl_matrix_set(A,i,q,c*aiq+s*aip); }
      for (int i=p+1;i<q;i++){
        double api=gsl_matrix_get(A,p,i);
        double aiq=gsl_matrix_get(A,i,q);
        gsl_matrix_set(A,p,i,c*api-s*aiq);
        gsl_matrix_set(A,i,q,c*aiq+s*api); }
      for (int i=q+1;i<n;i++){
        double api=gsl_matrix_get(A,p,i);
        double aqi=gsl_matrix_get(A,q,i);
        gsl_matrix_set(A,p,i,c*api-s*aqi);
        gsl_matrix_set(A,q,i,c*aqi+s*api); }
      for (int i=0;i<n;i++){
        double vip=gsl_matrix_get(V,i,p);
        double viq=gsl_matrix_get(V,i,q);
        gsl_matrix_set(V,i,p,c*vip-s*viq);
        gsl_matrix_set(V,i,q,c*viq+s*vip); }
    } } }while (changed!=0);
return sweeps; }

```



# Chapter 5

## Ordinary differential equations

### 5.1 Introduction

Ordinary differential equations (ODE) are generally defined as differential equations in one variable where the highest order derivative enters linearly. Such equations invariably arise in many different contexts throughout mathematics and science generally as soon as changes in the system at hand are considered, usually with respect to variations of certain parameters.

Ordinary differential equations can be generally reformulated as (coupled) systems of first-order ordinary differential equations,

$$\mathbf{y}'(x) = \mathbf{f}(x, \mathbf{y}) , \quad (5.1)$$

where  $\mathbf{y}' \doteq d\mathbf{y}/dx$ , and the variables  $\mathbf{y}$  and the *right-hand side function*  $\mathbf{f}(x, \mathbf{y})$  are understood as column-vectors. For example, a second order differential equation in the form

$$u'' = g(x, u, u') \quad (5.2)$$

can be rewritten as a system of two first-order equations,

$$\begin{cases} y_1' = y_2 \\ y_2' = g(x, y_1, y_2) \end{cases} , \quad (5.3)$$

using the variable substitution  $y_1 = u$ ,  $y_2 = u'$ .

In practice ODEs are usually supplemented with boundary conditions which pick out a certain class or a unique solution of the ODE. In the following we shall mostly consider the *initial value problem*: an ODE with the boundary condition in the form of an initial condition at a given point  $a$ ,

$$\mathbf{y}(a) = \mathbf{y}_0 . \quad (5.4)$$

The problem then is to find the value of the solution  $\mathbf{y}$  at some other point  $b$ .

Finding a solution to an ODE is often referred to as *integrating* the ODE.

An integration algorithm typically advances the solution from the initial point  $a$  to the final point  $b$  in a number of discrete steps

$$\{x_0 \doteq a, x_1, \dots, x_{n-1}, x_n \doteq b\}. \quad (5.5)$$

An efficient algorithm tries to integrate an ODE using as few steps as possible under the constraint of the given accuracy goal. For this purpose the algorithm should continuously adjust the step-size during the integration, using few larger steps in the regions where the solution is smooth and perhaps many smaller steps in more treacherous regions.

Typically, an adaptive step-size ODE integrator is implemented as two routines. One of them—called *driver*—monitors the local errors and tolerances and adjusts the step-sizes. To actually perform a step the driver calls a separate routine—the *stepper*—which advances the solution by one step, using one of the many available algorithms, and estimates the local error. The GNU Scientific Library, GSL, implements about a dozen of different steppers and a tunable adaptive driver.

In the following we shall discuss several of the popular driving algorithms and stepping methods for solving initial-value ODE problems.

## 5.2 Error estimate

In an adaptive step-size algorithm the stepping routine must provide an estimate of the integration error, upon which the driver bases its strategy to determine the optimal step-size for a user-specified accuracy goal.

A stepping method is generally characterized by its *order*: a method has order  $p$  if it can integrate exactly an ODE where the solution is a polynomial of order  $p$ . In other words, for small  $h$  the error of the order- $p$  method is  $O(h^{p+1})$ .

### 5.2.1 Runge's principle

For sufficiently small steps the error  $\delta y$  of an integration step for a method of a given order  $p$  can be estimated by comparing the solution  $\mathbf{y}_{\text{full\_step}}$ , obtained with one full-step integration, against a potentially more precise solution,  $\mathbf{y}_{\text{two\_half\_steps}}$ , obtained with two consecutive half-step integrations,

$$\delta \mathbf{y} = \frac{\mathbf{y}_{\text{full\_step}} - \mathbf{y}_{\text{two\_half\_steps}}}{2^p - 1}. \quad (5.6)$$

where  $p$  is the order of the algorithm used. Indeed, if the step-size  $h$  is small, we can assume

$$\delta \mathbf{y}_{\text{full\_step}} = Ch^{p+1}, \quad (5.7)$$

$$\delta \mathbf{y}_{\text{two\_half\_steps}} = 2C \left(\frac{h}{2}\right)^{p+1} = \frac{Ch^{p+1}}{2^p}, \quad (5.8)$$

where  $\delta \mathbf{y}_{\text{full\_step}}$  and  $\delta \mathbf{y}_{\text{two\_half\_steps}}$  are the errors of the full-step and two half-steps integrations, and  $C$  is an unknown constant. The two can be combined as

$$\begin{aligned} \mathbf{y}_{\text{full\_step}} - \mathbf{y}_{\text{two\_half\_steps}} &= \delta \mathbf{y}_{\text{full\_step}} - \delta \mathbf{y}_{\text{two\_half\_steps}} \\ &= \frac{Ch^{p+1}}{2^p} (2^p - 1), \end{aligned} \quad (5.9)$$

from which it follows that

$$\frac{Ch^{p+1}}{2^p} = \frac{\mathbf{y}_{\text{full\_step}} - \mathbf{y}_{\text{two\_half\_steps}}}{2^p - 1}. \quad (5.10)$$

One has, of course, to take the potentially more precise  $\mathbf{y}_{\text{two\_half\_steps}}$  as the approximation to the solution  $\mathbf{y}$ . Its error is then given as

$$\delta \mathbf{y}_{\text{two\_half\_steps}} = \frac{Ch^{p+1}}{2^p} = \frac{\mathbf{y}_{\text{full\_step}} - \mathbf{y}_{\text{two\_half\_steps}}}{2^p - 1}, \quad (5.11)$$

which had to be demonstrated. This prescription is often referred to as the *Runge's principle*.

One drawback of the Runge's principle is that the full-step and the two half-step calculations generally do not share evaluations of the right-hand side function  $\mathbf{f}(x, \mathbf{y})$ , and therefore many extra evaluations are needed to estimate the error.

### 5.2.2 Different orders

An alternative prescription for error estimation is to make the same step-size integration using two methods of *different orders*, with the difference between the two solutions providing the estimate of the error. If the lower order method mostly uses the same evaluations of the right-hand side function—in which case it is called *embedded* in the higher order method—the error estimate does not need additional evaluations.

Predictor-corrector methods are naturally of embedded type: the correction—which generally increases the order of the method—itself can serve as the estimate of the error.

### 5.3 Runge-Kutta methods

Runge-Kutta methods are one-step methods which advance the solution over the current step using only the information gathered from within the step itself. The solution  $\mathbf{y}$  is advanced from the point  $x_i$  to  $x_{i+1} = x_i + h$ , where  $h$  is the step-size, using a one-step formula,

$$\mathbf{y}_{i+1} = \mathbf{y}_i + h\mathbf{k}, \quad (5.12)$$

where  $\mathbf{y}_{i+1}$  is the approximation to  $\mathbf{y}(x_{i+1})$ , and the value  $\mathbf{k}$  is chosen such that the method integrates exactly an ODE whose solution is a polynomial of the highest possible order.

The Runge-Kutta methods are distinguished by their order. Again, a method has order  $p$  if it can integrate exactly an ODE where the solution is a polynomial of order  $p$  (or if for small  $h$  the error of the method is  $O(h^{p+1})$ ).

The first order Runge-Kutta method is the *Euler's method*,

$$\mathbf{k} = \mathbf{f}(x_0, \mathbf{y}_0). \quad (5.13)$$

Second order Runge-Kutta methods advance the solution by an auxiliary evaluation of the derivative. For example, the *mid-point method*,

$$\begin{aligned} \mathbf{k}_0 &= \mathbf{f}(x_0, \mathbf{y}_0), \\ \mathbf{k}_{1/2} &= \mathbf{f}(x_0 + \frac{1}{2}h, \mathbf{y}_0 + \frac{1}{2}h\mathbf{k}_0), \\ \mathbf{k} &= \mathbf{k}_{1/2}, \end{aligned} \quad (5.14)$$

or the *two-point method*, also called the *Heun's method*

$$\begin{aligned} \mathbf{k}_0 &= \mathbf{f}(x_0, \mathbf{y}_0), \\ \mathbf{k}_1 &= \mathbf{f}(x_0 + h, \mathbf{y}_0 + h\mathbf{k}_0), \\ \mathbf{k} &= \frac{1}{2}(\mathbf{k}_0 + \mathbf{k}_1). \end{aligned} \quad (5.15)$$

These two methods can be combined into a third order method,

$$\mathbf{k} = \frac{1}{6}\mathbf{k}_0 + \frac{4}{6}\mathbf{k}_{1/2} + \frac{1}{6}\mathbf{k}_1. \quad (5.16)$$

The most common is the fourth-order method, which is called *RK4* or simply *the Runge-Kutta method*,

$$\begin{aligned} \mathbf{k}_0 &= \mathbf{f}(x_0, \mathbf{y}_0), \\ \mathbf{k}_1 &= \mathbf{f}(x_0 + \frac{1}{2}h, \mathbf{y}_0 + \frac{1}{2}h\mathbf{k}_0), \\ \mathbf{k}_2 &= \mathbf{f}(x_0 + \frac{1}{2}h, \mathbf{y}_0 + \frac{1}{2}h\mathbf{k}_1), \\ \mathbf{k}_3 &= \mathbf{f}(x_0 + h, \mathbf{y}_0 + h\mathbf{k}_2), \\ \mathbf{k} &= \frac{1}{6}\mathbf{k}_0 + \frac{1}{3}\mathbf{k}_1 + \frac{1}{3}\mathbf{k}_2 + \frac{1}{6}\mathbf{k}_3. \end{aligned} \quad (5.17)$$

A general Runge-Kutta method can be written as

$$\mathbf{y}_{n+1} = \mathbf{y}_n + \sum_{i=1}^s b_i \mathbf{K}_i, \quad (5.18)$$

where

$$\begin{aligned} \mathbf{K}_1 &= h\mathbf{f}(x_n, \mathbf{y}_n), \\ \mathbf{K}_2 &= h\mathbf{f}(x_n + c_2h, \mathbf{y}_n + a_{21}\mathbf{K}_1), \\ \mathbf{K}_3 &= h\mathbf{f}(x_n + c_3h, \mathbf{y}_n + a_{31}\mathbf{K}_1 + a_{32}\mathbf{K}_2), \\ &\vdots \\ \mathbf{K}_s &= h\mathbf{f}(x_n + c_sh, \mathbf{y}_n + a_{s1}\mathbf{K}_1 + a_{s2}\mathbf{K}_2 + \cdots + a_{s,s-1}\mathbf{K}_{s-1}). \end{aligned} \quad (5.19)$$

Here the upper case  $\mathbf{K}_i$  are simply the lower case  $\mathbf{k}_i$  multiplied for convenience by the step-size  $h$ .

To specify a particular Runge-Kutta method one needs to provide the coefficients  $\{a_{ij} | 1 \leq j < i \leq s\}$ ,  $\{b_i | i = 1..s\}$  and  $\{c_i | i = 1..s\}$ . The matrix  $[a_{ij}]$  is called the Runge-Kutta matrix, while the coefficients  $b_i$  and  $c_i$  are known as the weights and the nodes. These data are usually arranged in the so called *Butcher's tableau*,

$$\begin{array}{c|cccc} 0 & & & & \\ c_2 & a_{21} & & & \\ c_3 & a_{31} & a_{32} & & \\ \vdots & \vdots & & \ddots & \\ c_s & a_{s1} & a_{s2} & \cdots & a_{s,s-1} \\ \hline & b_1 & b_2 & \cdots & b_{s-1} & b_s \end{array} \quad (5.20)$$

For example, the Butcher's tableau for the RK4 method is

$$\begin{array}{c|cccc} 0 & & & & \\ 1/2 & 1/2 & & & \\ 1/2 & 0 & 1/2 & & \\ 1 & 0 & 0 & 1 & \\ \hline & 1/6 & 1/3 & 1/3 & 1/6 \end{array} \quad (5.21)$$

### 5.3.1 Embedded methods with error estimates

The embedded Runge-Kutta methods in addition to advancing the solution by one step also produce an estimate of the local error of the step. This is done by having two methods in the tableau, one with a certain order  $p$  and another one with order  $p - 1$ .

Table 5.1: Embedded midpoint/Euler method with error estimate

```

void rkstep12(void f(int n, double x, double*yx, double*dydx),
int n, double x, double* yx, double h, double* yh, double* dy){
    int i; double k0[n], yt[n], k12[n]; /* VLA: gcc -std=c99 */
    f(n, x, yx, k0); for(i=0; i<n; i++) yt[i]=yx[i]+ k0[i]*h/2;
    f(n, x+h/2, yt, k12); for(i=0; i<n; i++) yh[i]=yx[i]+k12[i]*h;
    for(i=0; i<n; i++) dy[i]=(k0[i]-k12[i])*h/2; /* optimistic */
}

```

The difference between the two methods gives the the estimate of the local error. The lower order method is *embedded* in the higher order method, that is, it uses the same  $\mathbf{K}$ -values. This allows a very effective estimate of the error.

The embedded lower order method is written as

$$\mathbf{y}_{n+1}^* = \mathbf{y}_n + \sum_{i=1}^s b_i^* \mathbf{K}_i, \quad (5.22)$$

where  $\mathbf{K}_i$  are the same as for the higher order method. The error estimate is then given as

$$\mathbf{e}_n = \mathbf{y}_{n+1} - \mathbf{y}_{n+1}^* = \sum_{i=1}^s (b_i - b_i^*) \mathbf{K}_i. \quad (5.23)$$

The Butcher's tableau for this kind of method is extended by one row to give the values of  $b_i^*$ .

The simplest embedded methods are Heun-Euler method,

$$\begin{array}{c|cc} 0 & & \\ 1 & 1 & \\ \hline & 1/2 & 1/2 \\ & 1 & 0 \end{array}, \quad (5.24)$$

and midpoint-Euler method,

$$\begin{array}{c|cc} 0 & & \\ 1/2 & 1/2 & \\ \hline & 0 & 1 \\ & 1 & 0 \end{array}, \quad (5.25)$$

which both combine methods of orders 2 and 1. Table (5.1) shows a C-language implementation of the embedded midpoint/Euler method with error estimate.

The *Bogacki-Shampine method* [6] combines methods of orders 3 and 2,

$$\begin{array}{c|cccc}
 0 & & & & \\
 1/2 & 1/2 & & & \\
 3/4 & 0 & 3/4 & & \\
 1 & 2/9 & 1/3 & 4/9 & \\
 \hline
 & 2/9 & 1/3 & 4/9 & 0 \\
 & 7/24 & 1/4 & 1/3 & 1/8
 \end{array} . \tag{5.26}$$

Bogacki and Shampine argue that their method has better stability properties and actually outperforms higher order methods at lower accuracy goal calculations. This method has the FSAL—First Same As Last—property: the value  $\mathbf{k}_4$  at one step equals  $\mathbf{k}_1$  at the next step; thus only three function evaluations are needed per step. Following is a simple implementation which does not utilise this property for the sake of presentational clarity.

```

void rkstep23( void f(int n, double x, double* y, double* dydx),
int n, double x, double* yx, double h, double* yh, double* dy){
  int i; double k1[n], k2[n], k3[n], k4[n], yt[n]; /* VLA: -std=c99 */
  f(n, x, yx, k1); for(i=0; i<n; i++) yt[i]=yx[i]+1./2*k1[i]*h;
  f(n, x+1./2*h, yt, k2); for(i=0; i<n; i++) yt[i]=yx[i]+3./4*k2[i]*h;
  f(n, x+3./4*h, yt, k3); for(i=0; i<n; i++)
  yh[i]=yx[i]+(2./9 *k1[i]+1./3*k2[i]+4./9*k3[i])*h;
  f(n, x+h, yt, k4); for(i=0; i<n; i++){
  yt[i]=yx[i]+(7./24*k1[i]+1./4*k2[i]+1./3*k3[i]+1./8*k4[i])*h;
  dy[i]=yh[i]-yt[i];
  }
}

```

The Runge-Kutta-Fehlberg method [10]—called *RKF45*—implemented in the renowned `rkf45` Fortran routine, has two methods of orders 5 and 4:

$$\begin{array}{c|cccccc}
 0 & & & & & \\
 1/4 & 1/4 & & & & \\
 3/8 & 3/32 & 9/32 & & & \\
 12/13 & 1932/2197 & -7200/2197 & 7296/2197 & & \\
 1 & 439/216 & -8 & 3680/513 & -845/4104 & \\
 1/2 & -8/27 & 2 & -3544/2565 & 1859/4104 & -11/40 \\
 \hline
 & 16/135 & 0 & 6656/12825 & 28561/56430 & -9/50 & 2/55 \\
 & 25/216 & 0 & 1408/2565 & 2197/4104 & -1/5 & 0
 \end{array}$$

## 5.4 Multistep methods

Multistep methods try to use the information about the function gathered at the previous steps. They are generally not *self-starting* as there are no previous steps at the

start of the integration. The first step must be done with a one-step method like Runge-Kutta.

A number of multistep methods have been devised (and named after different mathematicians); we shall only consider a few simple ones here to get the idea of how it works.

### 5.4.1 Two-step method

Given the previous point,  $(x_{i-1}, \mathbf{y}_{i-1})$ , in addition to the current point  $(x_i, \mathbf{y}_i)$ , the sought function  $\mathbf{y}$  can be approximated in the vicinity of the point  $x_i$  as a second order polynomial,

$$\mathbf{y}(x) \approx \mathbf{p}_2(x) = \mathbf{y}_i + \mathbf{y}'_i \cdot (x - x_i) + \mathbf{c} \cdot (x - x_i)^2, \quad (5.27)$$

where  $\mathbf{y}'_i = \mathbf{f}(x_i, \mathbf{y}_i)$  and the coefficient  $\mathbf{c}$  can be found from the condition

$$\mathbf{p}_2(x_{i-1}) = \mathbf{y}_{i-1}, \quad (5.28)$$

which gives

$$\mathbf{c} = \frac{\mathbf{y}_{i-1} - \mathbf{y}_i + \mathbf{y}'_i \cdot (x_i - x_{i-1})}{(x_i - x_{i-1})^2}. \quad (5.29)$$

The value  $\mathbf{y}_{i+1}$  of the function at the next point,  $x_{i+1} \doteq x_i + h$ , can now be estimated as  $\mathbf{y}_{i+1} = \mathbf{p}_2(x_{i+1})$  from (5.27).

The error of this second-order two-step stepper can be estimated by a comparison with the first-order Euler's step, which is given by the linear part of (5.27). The correction term  $\mathbf{c}h^2$  can serve as the error estimate,

$$\delta\mathbf{y} = \mathbf{c}h^2. \quad (5.30)$$

### 5.4.2 Two-step method with extra evaluation

One can further increase the order of the approximation (5.27) by adding a third order term,

$$\mathbf{y}(x) \approx \mathbf{p}_3(x) = \mathbf{p}_2(x) + \mathbf{d} \cdot (x - x_i)^2(x - x_{i-1}). \quad (5.31)$$

The coefficient  $\mathbf{d}$  can be found from the matching condition at a certain point  $t$  inside the interval,

$$\mathbf{p}'_3(t) = \mathbf{f}(t, \mathbf{p}_2(t)) \doteq \mathbf{f}_2, \quad (5.32)$$

where  $x_i < t < x_i + h$ . This gives

$$\mathbf{d} = \frac{\mathbf{f}_2 - \mathbf{y}'_i - 2\mathbf{c} \cdot (t - x_i)}{2(t - x_i)(t - x_{i-1}) + (t - x_i)^2}. \quad (5.33)$$



The error estimate at the point  $x_{i+1} \doteq x_0 + h$  is again given as the difference between the higher and the lower order methods,

$$\delta \mathbf{y} = \mathbf{p}_3(x_{i+1}) - \mathbf{p}_2(x_{i+1}). \quad (5.34)$$

## 5.5 Predictor-corrector methods

A predictor-corrector method uses extra iterations to improve the solution. It is an algorithm that proceeds in two steps. First, the predictor step calculates a rough approximation of  $\mathbf{y}(x+h)$ . Second, the corrector step refines the initial approximation. Additionally the corrector step can be repeated in the hope that this achieves an even better approximation to the true solution.

For example, the two-point Runge-Kutta method (5.15) is as actually a predictor-corrector method, as it first calculates the *prediction*  $\tilde{\mathbf{y}}_{i+1}$  for  $\mathbf{y}(x_{i+1})$ ,

$$\tilde{\mathbf{y}}_{i+1} = \mathbf{y}_i + h\mathbf{f}(x_i, \mathbf{y}_i), \quad (5.35)$$

and then uses this prediction in a *correction* step,

$$\check{\mathbf{y}}_{i+1} = \mathbf{y}_i + h\frac{1}{2}(\mathbf{f}(x_i, \mathbf{y}_i) + \mathbf{f}(x_{i+1}, \tilde{\mathbf{y}}_{i+1})). \quad (5.36)$$

### 5.5.1 Two-step method with correction

Similarly, one can use the two-step approximation (5.27) as a predictor, and then improve it by one order with a correction step, namely

$$\check{\mathbf{y}}(x) = \bar{\mathbf{y}}(x) + \check{\mathbf{d}} \cdot (x - x_i)^2(x - x_{i-1}). \quad (5.37)$$

The coefficient  $\check{\mathbf{d}}$  can be found from the condition  $\check{\mathbf{y}}'(x_{i+1}) = \bar{\mathbf{f}}_{i+1}$ , where  $\bar{\mathbf{f}}_{i+1} \doteq \mathbf{f}(x_{i+1}, \bar{\mathbf{y}}(x_{i+1}))$ ,

$$\check{\mathbf{d}} = \frac{\bar{\mathbf{f}}_{i+1} - \mathbf{y}'_i - 2\mathbf{c} \cdot (x_{i+1} - x_i)}{2(x_{i+1} - x_i)(x_{i+1} - x_{i-1}) + (x_{i+1} - x_i)^2}. \quad (5.38)$$

Equation (5.37) gives a better estimate,  $\mathbf{y}_{i+1} = \check{\mathbf{y}}(x_{i+1})$ , of the sought function at the point  $x_{i+1}$ . In this context the formula (5.27) serves as *predictor*, and (5.37) as *corrector*. The difference between the two gives an estimate of the error.

This method is equivalent to the two-step method with an extra evaluation where the extra evaluation is done at the full step.

## 5.6 Adaptive step-size control

Let *tolerance*  $\tau$  be the maximal accepted error consistent with the required accuracy to be achieved in the integration of an ODE. Suppose the integration is done in  $n$  steps of size  $h_i$  such that  $\sum_{i=1}^n h_i = b - a$ . Under assumption that the errors at the integration steps are random and statistically uncorrelated, the local tolerance  $\tau_i$  for the step  $i$  has to scale as the square root of the step-size,

$$\tau_i = \tau \sqrt{\frac{h_i}{b-a}}. \quad (5.39)$$

Indeed, if the local error  $e_i$  on the step  $i$  is less than the local tolerance,  $e_i \leq \tau_i$ , the total error  $E$  will be consistent with the total tolerance  $\tau$ ,

$$E \approx \sqrt{\sum_{i=1}^n e_i^2} \leq \sqrt{\sum_{i=1}^n \tau_i^2} = \tau \sqrt{\sum_{i=1}^n \frac{h_i}{b-a}} = \tau. \quad (5.40)$$

The current step  $h_i$  is accepted if the local error  $e_i$  is smaller than the local tolerance  $\tau_i$ , after which the next step is attempted with the step-size adjusted according to the following empirical prescription [9],

$$h_{i+1} = h_i \times \left( \frac{\tau_i}{e_i} \right)^{\text{Power}} \times \text{Safety}, \quad (5.41)$$

where  $\text{Power} \approx 0.25$  and  $\text{Safety} \approx 0.95$ .

If the local error is larger than the local tolerance the step is rejected and a new step is attempted with the step-size adjusted according to the same prescription (5.41).

One simple prescription for the local tolerance  $\tau_i$  and the local error  $e_i$  to be used in (5.41) is

$$\tau_i = (\epsilon \|\mathbf{y}_i\| + \delta) \sqrt{\frac{h_i}{b-a}}, \quad e_i = \|\delta \mathbf{y}_i\|, \quad (5.42)$$

where  $\delta$  and  $\epsilon$  are the required absolute and relative precision and  $\delta \mathbf{y}_i$  is the estimate of the integration error at the step  $i$ .

A more elaborate prescription considers components of the solution separately,

$$(\tau_i)_k = (\epsilon |(\mathbf{y}_i)_k| + \delta) \sqrt{\frac{h_i}{b-a}}, \quad (\mathbf{e}_i)_k = |(\delta \mathbf{y}_i)_k|, \quad (5.43)$$

where the index  $k$  runs over the components of the solution. In this case the step acceptance criterion also becomes component-wise: the step is accepted, if

$$\forall k : (\mathbf{e}_i)_k < (\tau_i)_k. \quad (5.44)$$

The factor  $\tau_i/e_i$  in the step adjustment formula (5.41) is then replaced by

$$\frac{\tau_i}{e_i} \rightarrow \min_k \frac{(\tau_i)_k}{(\mathbf{e}_i)_k}. \quad (5.45)$$

Yet another refinement is to include the derivatives  $\mathbf{y}'$  of the solution into the local tolerance estimate, either overall,

$$\tau_i = \left( \epsilon\alpha \|\mathbf{y}_i\| + \epsilon\beta \|\mathbf{y}'_i\| + \delta \right) \sqrt{\frac{h_i}{b-a}}, \quad (5.46)$$

or component-wise,

$$(\tau_i)_k = \left( \epsilon\alpha |(\mathbf{y}_i)_k| + \epsilon\beta |(\mathbf{y}'_i)_k| + \delta \right) \sqrt{\frac{h_i}{b-a}}. \quad (5.47)$$

The weights  $\alpha$  and  $\beta$  are chosen by the user.

Following is a simple C-language implementation of the described algorithm.

```

int ode_driver(void f(int n, float x, float*y, float*dydx),
int n, float*xlist, float**ylist,
float b, float h, float acc, float eps, int max){
  int i, k=0; float x,*y,s, err, normy, tol, a=xlist[0], yh[n], dy[n];
  while( xlist[k]<b){
    x=xlist[k], y=ylist[k]; if(x+h>b) h=b-x;
    ode_stepper(f,n,x,y,h,yh,dy);
    s=0; for(i=0;i<n;i++) s+=dy[i]*dy[i]; err =sqrt(s);
    s=0; for(i=0;i<n;i++) s+=yh[i]*yh[i]; normy=sqrt(s);
    tol=(normy*eps+acc)*sqrt(h/(b-a));
    if(err<tol){ /* accept step and continue */
      k++; if(k>max-1) return -k; /* uups */
      xlist[k]=x+h; for(i=0;i<n;i++)ylist[k][i]=yh[i];
    }
    if(err>0) h*=pow(tol/err,0.25)*0.95; else h*=2;
  } /* end while */
  return k+1; } /* return the number of entries in xlist/ylist */

```



# Chapter 6

## Nonlinear equations

### 6.1 Introduction

*Non-linear equations* (or *root-finding*) is a problem of finding a set of  $n$  variables  $\mathbf{x} = \{x_1, \dots, x_n\}$  which satisfy a system of  $n$  non-linear equations

$$f_i(x_1, \dots, x_n) = 0 \Big|_{i=1, \dots, n} . \quad (6.1)$$

In matrix notation the system is written as

$$\mathbf{f}(\mathbf{x}) = 0 , \quad (6.2)$$

where  $\mathbf{f}(\mathbf{x}) \doteq \{f_1(x_1, \dots, x_n), \dots, f_n(x_1, \dots, x_n)\}$ .

In one-dimension,  $n = 1$ , it is generally possible to plot the function in the region of interest and see whether the graph crosses the  $x$ -axis. One can then be sure the root exists and even figure out its approximate position to start one's root-finding algorithm from. In multi-dimensions one generally does not know if the root exists at all, until it is found.

The root-finding algorithms generally proceed by iteration, starting from some approximate solution and making consecutive steps—hopefully in the direction of the suspected root—until some convergence criterion is satisfied. The procedure is generally not even guaranteed to converge unless starting from a point close enough to the sought root.

We shall only consider the multi-dimensional case here since i) the multi-dimensional root-finding is more difficult, and ii) the multi-dimensional routines can also be used in the one-dimensional case.

## 6.2 Newton's method

Newton's method (also referred to as Newton-Raphson method, after Isaac Newton and Joseph Raphson) is a root-finding algorithm that uses the first term of the Taylor series of the functions  $f_i$  to linearise the system (6.1) in the vicinity of a suspected root. It is one of the oldest and best known methods and is a basis of a number of more refined methods.

Suppose that the point  $\mathbf{x} = \{x_1, \dots, x_n\}$  is close to the root. The Newton's algorithm tries to find the step  $\Delta\mathbf{x}$  which would move the point towards the root, such that

$$f_i(\mathbf{x} + \Delta\mathbf{x}) = 0 \Big|_{i=1, \dots, n}. \quad (6.3)$$

The first order Taylor expansion of (6.3) gives a system of linear equations,

$$f_i(\mathbf{x}) + \sum_{k=1}^n \frac{\partial f_i}{\partial x_k} \Delta x_k = 0 \Big|_{i=1, \dots, n}, \quad (6.4)$$

or, in the matrix form,

$$\mathbf{J}\Delta\mathbf{x} = -\mathbf{f}(\mathbf{x}), \quad (6.5)$$

where  $\mathbf{J}$  is the matrix of partial derivatives,

$$J_{ik} \doteq \frac{\partial f_i}{\partial x_k}, \quad (6.6)$$

called the *Jacobian matrix*. In practice, if derivatives are not available analytically, one uses finite differences,

$$\frac{\partial f_i}{\partial x_k} \approx \frac{f_i(x_1, \dots, x_{k-1}, x_k + \delta x, x_{k+1}, \dots, x_n) - f_i(x_1, \dots, x_k, \dots, x_n)}{\delta x}, \quad (6.7)$$

with  $\delta x \ll s$  with  $s$  being the typical scale of the problem at hand.

One should always rescale one's problem such that the typical scale is around unity. In this case  $\delta x$  can be chosen as the square root of the machine epsilon.

The solution  $\Delta\mathbf{x}$  to the linear system (6.5)—called the Newton's step—gives the approximate direction and the approximate step-size towards the solution.

The Newton's method converges quadratically if sufficiently close to the solution. Otherwise the full Newton's step  $\Delta\mathbf{x}$  might actually diverge from the solution. Therefore in practice a more conservative step,  $\lambda\Delta\mathbf{x}$  with  $\lambda < 1$ , is usually taken. The strategy of finding the optimal  $\lambda$  is referred to as *line search*.

It is typically not worth the effort to find  $\lambda$  which minimizes  $\|\mathbf{f}(\mathbf{x} + \lambda\Delta\mathbf{x})\|$  exactly, since  $\Delta\mathbf{x}$  is only an approximate direction towards the root. Instead an inexact but

quick minimization strategy is usually used, like the *backtracking line search* where one first attempts the full step,  $\lambda = 1$ , and then backtracks,  $\lambda \leftarrow \lambda/2$ , until the condition

$$\|\mathbf{f}(\mathbf{x} + \lambda\Delta\mathbf{x})\| < \left(1 - \frac{\lambda}{2}\right) \|\mathbf{f}(\mathbf{x})\| \quad (6.8)$$

is satisfied. If the condition is not satisfied for sufficiently small  $\lambda_{\min}$  the step is taken with  $\lambda_{\min}$  simply to step away from the difficult place and try again.

Following is a typical algorithm of the Newton's method with backtracking line search and condition (6.8),

```

repeat
  calculate the Jacobian matrix J
  solve J $\Delta\mathbf{x} = -\mathbf{f}(\mathbf{x})$  for  $\Delta\mathbf{x}$ 
   $\lambda \leftarrow 1$ 
  while (  $\|\mathbf{f}(\mathbf{x} + \lambda\Delta\mathbf{x})\| > (1 - \frac{\lambda}{2}) \|\mathbf{f}(\mathbf{x})\|$  and  $\lambda > \frac{1}{64}$  ) do  $\lambda \leftarrow \lambda/2$ 
   $\mathbf{x} \leftarrow \mathbf{x} + \lambda\Delta\mathbf{x}$ 
until converged (e.g.  $\|\mathbf{f}(\mathbf{x})\| < \text{tolerance}$ )

```

A somewhat more refined backtracking linesearch is based on an approximate minimization of the function

$$g(\lambda) \doteq \frac{1}{2} \|\mathbf{f}(\mathbf{x} + \lambda\Delta\mathbf{x})\|^2 \quad (6.9)$$

using interpolation. The values  $g(0) = \frac{1}{2} \|\mathbf{f}(\mathbf{x})\|^2$  and  $g'(0) = -\|\mathbf{f}(\mathbf{x})\|^2$  are already known (check this). If the previous step with certain  $\lambda_{\text{trial}}$  was rejected, we also have  $g(\lambda_{\text{trial}})$ . These three quantities allow to build a quadratic approximation,

$$g(\lambda) \approx g(0) + g'(0)\lambda + c\lambda^2, \quad (6.10)$$

where

$$c = \frac{g(\lambda_{\text{trial}}) - g(0) - g'(0)\lambda_{\text{trial}}}{\lambda_{\text{trial}}^2}. \quad (6.11)$$

The minimum of this approximation (determined by the condition  $g'(\lambda) = 0$ ),

$$\lambda_{\text{next}} = -\frac{g'(0)}{2c}, \quad (6.12)$$

becomes the next trial step-size.

The procedure is repeated recursively until either condition (6.8) is satisfied or the step becomes too small (in which case it is taken unconditionally in order to simply get away from the difficult place).

## 6.3 Quasi-Newton methods

The Newton's method requires calculation of the Jacobian matrix at every iteration. This is generally an expensive operation. Quasi-Newton methods avoid calculation of the Jacobian matrix at the new point  $\mathbf{x} + \Delta\mathbf{x}$ , instead trying to use certain approximations, typically rank-1 updates.

### 6.3.1 Broyden's method

Broyden's algorithm [7] estimates the Jacobian  $\mathbf{J} + \Delta\mathbf{J}$  at the point  $\mathbf{x} + \Delta\mathbf{x}$  using the finite-difference approximation,

$$(\mathbf{J} + \Delta\mathbf{J})\Delta\mathbf{x} = \Delta\mathbf{f}, \quad (6.13)$$

where  $\Delta\mathbf{f} \doteq \mathbf{f}(\mathbf{x} + \Delta\mathbf{x}) - \mathbf{f}(\mathbf{x})$  and  $\mathbf{J}$  is the Jacobian at the point  $\mathbf{x}$ .

The matrix equation (6.13) is under-determined in more than one dimension as it contains only  $n$  equations to determine  $n^2$  matrix elements of  $\Delta\mathbf{J}$ . Broyden suggested to choose  $\Delta\mathbf{J}$  as a rank-1 update, linear in  $\Delta\mathbf{x}$ ,

$$\Delta\mathbf{J} = \mathbf{c} \Delta\mathbf{x}^T, \quad (6.14)$$

where the unknown vector  $\mathbf{c}$  can be found by substituting (6.14) into (6.13), which gives

$$\Delta\mathbf{J} = \frac{(\Delta\mathbf{f} - \mathbf{J}\Delta\mathbf{x})\Delta\mathbf{x}^T}{\Delta\mathbf{x}^T\Delta\mathbf{x}}. \quad (6.15)$$

In practice if one wanders too far from the point where  $\mathbf{J}$  was first calculated the accuracy of the updates may decrease significantly. In such case one might need to recalculate  $\mathbf{J}$  anew. For example, two successive steps with  $\lambda_{\min}$  might be interpreted as a sign of accuracy loss in  $\mathbf{J}$  and subsequently trigger its recalculation.

It also possible to update directly the inverse  $\mathbf{B} \doteq \mathbf{J}^{-1}$  of the Jacobian matrix using the Sherman-Morrison formula for the inverse of a rank-1 updated matrix,

$$(A + uv^T)^{-1} = A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u}. \quad (6.16)$$

For the update (6.15) the Sherman-Morrison formula gives

$$\Delta\mathbf{B} = \frac{\Delta\mathbf{x} - \mathbf{B}\Delta\mathbf{f}}{\Delta\mathbf{x}^T\mathbf{B}\Delta\mathbf{f}}\Delta\mathbf{x}^T\mathbf{B}. \quad (6.17)$$

Another update from the Broyden's class is

$$\Delta\mathbf{B} = \frac{\Delta\mathbf{x} - \mathbf{B}\Delta\mathbf{f}}{\Delta\mathbf{x}^T\Delta\mathbf{f}}\Delta\mathbf{x}^T. \quad (6.18)$$



Table 6.1: Python implementation of Newton's root-finding algorithm with backtracking.

```
def newton(f:"function",xstart:vector,eps:float=1e-3,dx:float=1e-6):
    x=xstart.copy(); n=x.size; J = matrix(n,n)
    while True :
        fx=f(x)
        for j in range(n) :
            x[j]+=dx
            df=f(x)-fx
            for i in range(n) : J[i,j] = df[i]/dx
            x[j]-=dx
        givens.qr(J)
        Dx = givens.solve(J,-fx)
        s=2
        while True :
            s/=2
            y=x+Dx*s
            fy=f(y)
            if fy.norm()<(1-s/2)*fx.norm() or s<0.02 : break
        x=y; fx=fy
        if Dx.norm()<dx or fx.norm()<eps : break
    return x;
```



# Chapter 7

## Numerical integration

### 7.1 Introduction

The term *numerical integration* refers to a broad family of algorithms to compute a numerical approximation to a definite (Riemann) integral.

Generally, the integral is approximated by a weighted sum of function values within the domain of integration,

$$\int_a^b f(x)dx \approx \sum_{i=1}^n w_i f(x_i). \quad (7.1)$$

Expression (7.1) is often referred to as *quadrature* (*cubature* for multidimensional integrals) or *rule*. The abscissas  $x_i$  (also called *nodes*) and the weights  $w_i$  of a quadrature are usually optimized—using one of a large number of different strategies—to suit a particular class of integration problems.

The best quadrature algorithm for a given problem depends on several factors, in particular on the integrand. Different classes of integrands generally require different quadratures for the most effective calculation.

A popular numerical integration library is QUADPACK [21]. It includes general purpose routines—like QAGS, based on an adaptive Gauss–Kronrod quadrature with acceleration—as well as a number of specialized routines. The GNU scientific library [9] (GSL) implements most of the QUADPACK routines and in addition includes a modern general-purpose adaptive routine CQUAD based on Clenshaw–Curtis quadratures [14].

In the following we shall consider some of the popular numerical integration algorithms.

## 7.2 Rectangle and trapezium rules

In mathematics, the *Riemann integral* is generally defined in terms of *Riemann sums* [22]. If the integration interval  $[a, b]$  is partitioned into  $n$  subintervals,

$$a = t_0 < t_1 < t_2 < \cdots < t_n = b. \quad (7.2)$$

the Riemann sum is defined as

$$\sum_{i=1}^n f(x_i) \Delta x_i, \quad (7.3)$$

where  $x_i \in [t_{i-1}, t_i]$  and  $\Delta x_i = t_i - t_{i-1}$ . Geometrically a Riemann sum can be interpreted as the area of a collection of adjacent rectangles with widths  $\Delta x_i$  and heights  $f(x_i)$ .

The Riemann integral is defined as the limit of a Riemann sum as the *mesh*—the length of the largest subinterval—of the partition approaches zero. Specifically, the number denoted as

$$\int_a^b f(x) dx \quad (7.4)$$

is called the Riemann integral, if for any  $\epsilon > 0$  there exists  $\delta > 0$  such that for any partition (7.2) with  $\max \Delta x_i < \delta$  we have

$$\left| \sum_{i=1}^n f(x_i) \Delta x_i - \int_a^b f(x) dx \right| < \epsilon. \quad (7.5)$$

A definite integral can be interpreted as the net signed area bounded by the graph of the integrand.

Now, the  $n$ -point *rectangle quadrature* is simply the Riemann sum (7.3),

$$\int_a^b f(x) dx \approx \sum_{i=1}^n f(x_i) \Delta x_i, \quad (7.6)$$

where the node  $x_i$  is often (but not always) taken in the middle of the corresponding subinterval,  $x_i = t_{i-1} + \frac{1}{2} \Delta x_i$ , and the subintervals are often (but not always) chosen equal,  $\Delta x_i = (b - a)/n$ . Geometrically the  $n$ -point rectangle rule is an approximation to the integral given by the area of a collection of  $n$  adjacent equal rectangles whose heights are determined by the values of the function (at the middle of the rectangle).

An  $n$ -point *trapezium rule* uses instead a collection of trapezia fitted under the graph,

$$\int_a^b f(x) dx \approx \sum_{i=1}^n \frac{f(t_{i-1}) + f(t_i)}{2} \Delta x_i. \quad (7.7)$$

Importantly, the trapezium rule is the average of two Riemann sums,

$$\sum_{i=1}^n \frac{f(t_{i-1}) + f(t_i)}{2} \Delta x_i = \frac{1}{2} \sum_{i=1}^n f(t_{i-1}) \Delta x_i + \frac{1}{2} \sum_{i=1}^n f(t_i) \Delta x_i. \quad (7.8)$$

Rectangle and trapezium quadratures both have the important feature of closely following the very mathematical definition of the integral as the limit of the Riemann sums. Therefore—disregarding the round-off errors—these two rules cannot fail if the integral exists.

For certain partitions of the interval the rectangle and trapezium rules coincide. For example, for the nodes

$$x_i = a + (b - a) \frac{i - \frac{1}{2}}{n}, \quad i = 1, \dots, n \quad (7.9)$$

both rules give the same quadrature with equal weights,  $w_i = (b - a)/n$ ,

$$\int_a^b f(x) dx \approx \frac{b - a}{n} \sum_{i=1}^n f\left(a + (b - a) \frac{i - \frac{1}{2}}{n}\right). \quad (7.10)$$

Rectangle and trapezium quadratures are rarely used on their own—because of the slow convergence—but they often serve as the basis for more advanced quadratures, for example adaptive quadratures and variable transformation quadratures considered below.

### 7.3 Quadratures with regularly spaced abscissas

A quadrature (7.1) with  $n$  predefined nodes  $x_i$  has  $n$  free parameters: the weights  $w_i$ . A set of  $n$  parameters can generally be tuned to satisfy  $n$  conditions. The archetypal set of conditions in quadratures is that the quadrature integrates exactly a set of  $n$  functions,

$$\{\phi_1(x), \dots, \phi_n(x)\}. \quad (7.11)$$

This leads to a set of  $n$  equations,

$$\sum_{i=1}^n w_i \phi_k(x_i) = I_k \Big|_{k=1, \dots, n}, \quad (7.12)$$

where the integrals

$$I_k \doteq \int_a^b \phi_k(x) dx \quad (7.13)$$

are assumed to be known. Equations (7.12) are linear in  $w_i$  and can be easily solved.

Since integration is a linear operation, the quadrature will then also integrate exactly any linear combination of functions (7.11).

A popular choice for predefined nodes is a *closed set*—that is, including the endpoints of the interval—of evenly spaced abscissas,

$$x_i = a + \frac{i-1}{n-1}(b-a) \Big|_{i=1, \dots, n} . \quad (7.14)$$

However, in practice it often happens that the integrand has an integrable singularity at one or both ends of the interval. In this case one can choose an *open set* of equidistant nodes,

$$x_i = a + \frac{i-\frac{1}{2}}{n}(b-a) \Big|_{i=1, \dots, n} . \quad (7.15)$$

The set of functions to be integrated exactly is generally chosen to suite the properties of the integrands at hand: the integrands must be well represented by linear combinations of the chosen functions.

### 7.3.1 Classical quadratures

Suppose the integrand can be well represented by the first few terms of its Taylor series,

$$f(x) = \sum_{k=0}^{\infty} \frac{f^{(k)}(a)}{k!} (x-a)^k , \quad (7.16)$$

where  $f^{(k)}$  is the  $k$ -th derivative of the integrand. This is often the case for analytic—that is, infinitely differentiable—functions. For such integrands one can obviously choose polynomials

$$\{1, x, x^2, \dots, x^{n-1}\} \quad (7.17)$$

as the set of functions to be integrated exactly.

This leads to the so called *classical quadratures*: quadratures with regularly spaced abscissas and *polynomials* as exactly integrable functions.

An  $n$ -point classical quadrature integrates exactly the first  $n$  terms of the function's Taylor expansion (7.16). The  $x^n$  order term will not be integrated exactly and will lead to an error of the quadrature. Thus the error  $E_n$  of the  $n$ -point classical quadrature is on the order of the integral of the  $x^n$  term in (7.16),

$$E_n \approx \int_a^b \frac{f^{(n)}(a)}{n!} (x-a)^n dx = \frac{f^{(n)}(a)}{(n+1)!} h^{n+1} \propto h^{n+1} , \quad (7.18)$$

where  $h = b - a$  is the length of the integration interval. A quadrature with the error of the order  $h^{n+1}$  is often called a *degree- $n$*  quadrature.

Table 7.1: Maxima script to calculate analytically the weights of an  $n$ -point classical quadrature with predefined abscissas in the interval  $[0, 1]$ .

```
n: 8; xs: makelist((i-1)/(n-1),i,1,n); /* nodes: adapt to your needs */
ws: makelist(concat(w,i),i,1,n);
ps: makelist(x^i,i,0,n-1); /* polynomials */
fs: makelist(buildq([i:i,ps:ps],lambda([x],ps[i])),i,1,n);
integ01: lambda([f],integrate(f(x),x,0,1));
Is: maplist(integ01,fs); /* calculate the integrals */
eq: lambda([f],lreduce("+",maplist(f,xs)*ws));
eqs: maplist(eq,fs)-Is; /* build equations */
solve(eqs,ws); /* solve for the weights */
```

---

If the integrand is smooth enough and the length  $h$  is small enough a classical quadrature with not so large  $n$  can provide a good approximation for the integral. However, for large  $n$  the weights of classical quadratures tend to have alternating signs, which leads to large round-off errors, which in turn negates the potentially higher accuracy of the quadrature. Again, if the integrand violates the assumption of Taylor expansion—for example by having an integrable singularity inside the integration interval—the higher order quadratures may perform poorly.

Classical quadratures are mostly of historical interest nowadays. Alternative methods—such as quadratures with optimized abscissas, adaptive, and variable transformation quadratures—are more stable and accurate and are normally preferred to classical quadratures.

Classical quadratures with equally spaced abscissas—both closed and open sets—are generally referred to as *Newton-Cotes quadratures*. An interested reader can generate Newton—Cotes quadratures of any degree  $n$  using the Maxima script in Table (7.1).

## 7.4 Quadratures with optimized abscissas

In quadratures with optimized abscissas not only the weights  $w_i$  but also the abscissas  $x_i$  are chosen optimally. The number of free parameters is thus  $2n$  and one can choose a set of  $2n$  functions,

$$\{\phi_1(x), \dots, \phi_{2n}(x)\}, \quad (7.19)$$

to be integrated exactly. This gives a system of  $2n$  equations, linear in  $w_i$  and non-linear in  $x_i$ ,

$$\sum_{i=1}^n w_i \phi_k(x_i) = I_k \Big|_{k=1, \dots, 2n}, \quad (7.20)$$

where again

$$I_k \doteq \int_a^b \phi_k(x) dx . \quad (7.21)$$

The weights and abscissas of the quadrature can be determined by solving this system of equations<sup>1</sup>.

Although quadratures with optimized abscissas are generally of much higher order,  $2n - 1$  compared to  $n - 1$  for non-optimal abscissas, the optimal points generally can not be reused at the next iteration in an adaptive algorithm.

### 7.4.1 Gauss quadratures

Gauss quadratures deal with a slightly more general form of integrals,

$$\int_a^b \omega(x) f(x) dx , \quad (7.23)$$

where  $\omega(x)$  is a positive weight function. For  $\omega(x) = 1$  the problem is the same as considered above. Popular choices of the weight function include  $\omega(x) = (1 - x^2)^{\pm 1/2}$ ,  $\exp(-x)$ ,  $\exp(-x^2)$  and others. The idea is to represent the integrand as a product  $\omega(x)f(x)$  such that all the difficulties go into the weight function  $\omega(x)$  while the remaining factor  $f(x)$  is smooth and well represented by polynomials.

An  $N$ -point *Gauss quadrature* is a quadrature with optimized abscissas,

$$\int_a^b \omega(x) f(x) dx \approx \sum_{i=1}^N w_i f(x_i) , \quad (7.24)$$

which integrates exactly a set of  $2N$  polynomials of the orders  $1, \dots, 2N - 1$  with the given weight  $\omega(x)$ .

#### Fundamental theorem

There is a theorem stating that there exists a set of polynomials  $p_n(x)$ , orthogonal on the interval  $[a, b]$  with the weight function  $\omega(x)$ ,

$$\int_a^b \omega(x) p_n(x) p_k(x) \propto \delta_{nk} . \quad (7.25)$$

---

<sup>1</sup>Here is, for example, an  $n = 2$  quadrature with optimized abscissas,

$$\int_{-1}^1 f(x) dx \approx f\left(-\sqrt{\frac{1}{3}}\right) + f\left(+\sqrt{\frac{1}{3}}\right) . \quad (7.22)$$



Now, one can prove that the optimal nodes for the  $N$ -point Gauss quadrature are the roots of the polynomial  $p_N(x)$ ,

$$p_N(x_i) = 0. \quad (7.26)$$

The idea behind the proof is to consider the integral

$$\int_a^b \omega(x)q(x)p_N(x)dx = 0, \quad (7.27)$$

where  $q(x)$  is an arbitrary polynomial of degree less than  $N$ . The quadrature should represent this integral exactly,

$$\sum_{i=1}^N w_i q(x_i) p_N(x_i) = 0. \quad (7.28)$$

Apparently this is only possible if  $x_i$  are the roots of  $p_N$ .

### Calculation of nodes and weights

A neat algorithm—usually referred to as Golub-Welsch algorithm [13]—for calculation of the nodes and weights of a Gauss quadrature is based on the symmetric form of the three-term recurrence relation for orthogonal polynomials,

$$xp_{n-1}(x) = \beta_n p_n(x) + \alpha_n p_{n-1}(x) + \beta_{n-1} p_{n-2}(x), \quad (7.29)$$

where  $p_{-1}(x) \doteq 0$ ,  $p_1(x) \doteq 1$ , and  $n = 1, \dots, N$ . This recurrence relation can be written in the matrix form,

$$x\mathbf{p}(x) = \mathbf{J}\mathbf{p}(x) + \beta_N p_N(x)\mathbf{e}_N, \quad (7.30)$$

where  $\mathbf{p}(x) \doteq \{p_0(x), \dots, p_{N-1}(x)\}^T$ ,  $\mathbf{e}_N = \{0, \dots, 0, 1\}^T$ , and the tridiagonal matrix  $\mathbf{J}$  — usually referred to as *Jacobi matrix* or *Jacobi operator* — is given as

$$\mathbf{J} = \begin{pmatrix} \alpha_1 & \beta_1 & & & \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \beta_2 & \alpha_3 & \beta_3 & \\ & & \ddots & \ddots & \\ & & & \beta_{N-1} & \alpha_N \end{pmatrix}. \quad (7.31)$$

Substituting the roots  $x_i$  of  $p_N$  — that is, the set  $\{x_i \mid p_N(x_i) = 0\}$  — into the matrix equation (7.30) leads to eigenvalue problem for the Jacobi matrix,

$$\mathbf{J}\mathbf{p}(x_i) = x_i \mathbf{p}(x_i). \quad (7.32)$$

Thus, the nodes of an  $N$ -point Gauss quadrature (the roots of the polynomial  $p_N$ ) are the eigenvalues of the Jacobi matrix  $J$  and can be calculated by a standard diagonalization<sup>2</sup> routine.

The weights can be obtained considering  $N$  integrals,

$$\int_a^b \omega(x)p_n(x)dx = \delta_{n0} \int_a^b \omega(x)dx, \quad n = 0, \dots, N-1. \quad (7.33)$$

Applying our quadrature gives the matrix equation,

$$\mathbf{P}\mathbf{w} = \mathbf{e}_1 \int_a^b \omega(x)dx, \quad (7.34)$$

where  $\mathbf{w} \doteq \{w_1, \dots, w_N\}^T$ ,  $\mathbf{e}_1 = \{1, 0, \dots, 0\}^T$ , and

$$\mathbf{P} \doteq \begin{pmatrix} p_0(x_1) & \dots & p_0(x_N) \\ p_1(x_1) & \dots & p_1(x_N) \\ \dots & \dots & \dots \\ p_{N-1}(x_1) & \dots & p_{N-1}(x_N) \end{pmatrix}. \quad (7.35)$$

Equation (7.34) is linear in  $w_i$  and can be solved directly. However, if diagonalization of the Jacobi matrix provided the normalized eigenvectors, the weights can be readily obtained using the following method.

The matrix  $\mathbf{P}$  apparently consists of non-normalized column eigenvectors of the matrix  $\mathbf{J}$ . The eigenvectors are orthogonal and therefore  $\mathbf{P}^T\mathbf{P}$  is a diagonal matrix with positive elements. Multiplying (7.34) by  $\mathbf{P}^T$  and then by  $(\mathbf{P}^T\mathbf{P})^{-1}$  from the left gives

$$\mathbf{w} = (\mathbf{P}^T\mathbf{P})^{-1}\mathbf{P}^T\mathbf{e}_1 \int_a^b \omega(x)dx. \quad (7.36)$$

From  $p_0(x) = 1$  it follows that  $\mathbf{P}^T\mathbf{e}_1 = \{1, \dots, 1\}^T$  and therefore

$$w_i = \frac{1}{(\mathbf{P}^T\mathbf{P})_{ii}} \int_a^b \omega(x)dx. \quad (7.37)$$

Let the matrix  $\mathbf{V}$  be the set of the normalized column eigenvectors of the matrix  $\mathbf{J}$ . The matrix  $\mathbf{V}$  is then connected with the matrix  $\mathbf{P}$  through the normalization equation,

$$\mathbf{V} = \sqrt{(\mathbf{P}^T\mathbf{P})^{-1}}\mathbf{P}. \quad (7.38)$$

Therefore, again taking into account that  $p_0(x) = 1$ , equation (7.37) can be written as

$$w_i = (V_{1i})^2 \int_a^b \omega(x)dx. \quad (7.39)$$

---

<sup>2</sup>A symmetric tridiagonal matrix can be diagonalized very effectively using the QR/RL algorithm.

Table 7.2: An Octave function that calculates the nodes and weights of the  $N$ -point Gauss-Legendre quadrature and then integrates a given function.

```

function Q = gauss_legendre(f,a,b,N)
beta = .5./sqrt(1-(2*(1:N-1)).^(-2)); % recurrence relation
J = diag(beta,1) + diag(beta,-1); % Jacobi matrix
[V,D] = eig(J); % diagonalization of J
x = diag(D); [x,i] = sort(x); % sorted nodes
w = V(1,i).^2*2; % weights
Q = w*f((a+b)/2+(b-a)/2**x)*(b-a)/2; % integral
endfunction;

```

### Example: Gauss-Legendre quadrature

Gauss-Legendre quadrature deals with the weight  $\omega(x) = 1$  on the interval  $[-1, 1]$ . The associated polynomials are Legendre polynomials  $\mathcal{P}_n(x)$ , hence the name. Their recurrence relation is usually given as

$$(2n - 1)x\mathcal{P}_{n-1}(x) = n\mathcal{P}_n(x) + (n - 1)\mathcal{P}_{n-2}(x). \quad (7.40)$$

Rescaling the polynomials (preserving  $p_0(x) = 1$ ) as

$$\sqrt{2n + 1}\mathcal{P}_n(x) = p_n(x) \quad (7.41)$$

reduces this recurrence relation to the symmetric form (7.29),

$$xp_{n-1}(x) = \frac{1}{2} \frac{1}{\sqrt{1 - (2n)^{-2}}} p_n(x) + \frac{1}{2} \frac{1}{\sqrt{1 - (2(n-1))^{-2}}} p_{n-2}(x). \quad (7.42)$$

Correspondingly, the coefficients in the matrix  $\mathbf{J}$  are

$$\alpha_n = 0, \quad \beta_n = \frac{1}{2} \frac{1}{\sqrt{1 - (2n)^{-2}}}. \quad (7.43)$$

The problem of finding the nodes and the weights of the  $N$ -point Gauss-Legendre quadrature is thus reduced to the eigenvalue problem for the Jacobi matrix with coefficients (7.43).

As an illustration of this algorithm Table (7.2) shows an Octave function which calculates the nodes and the weights of the  $N$ -point Gauss-Legendre quadrature and then integrates a given function.

### 7.4.2 Gauss-Kronrod quadratures

Generally, the error of a numerical integration is estimated by comparing the results from two rules of different orders. However, for ordinary Gauss quadratures the nodes for two rules of different orders almost never coincide. This means that one can not reuse the points of the lower order rule when calculating the higher order rule.

*Gauss-Kronrod* algorithm [18] remedies this inefficiency. The points inherited from the lower order rule are reused in the higher order rule as predefined nodes (with  $n$  weights as free parameters), and then  $m$  more optimal points are added ( $m$  abscissas and  $m$  weights as free parameters). The order of the method is  $n + 2m - 1$ . The lower order rule becomes *embedded*—that is, it uses a subset of the nodes—into the higher order rule. On the next iteration the procedure is repeated.

Patterson [20] has tabulated nodes and weights for several sequences of embedded Gauss-Kronrod rules.

## 7.5 Adaptive quadratures

Higher order quadratures suffer from round-off errors as the weights  $w_i$  generally have alternating signs. Again, using high order polynomials is dangerous as they typically oscillate wildly and may lead to Runge's phenomenon. Therefore, if the error of the quadrature is yet too large for a quadrature with sufficiently large  $n$ , the best strategy is to subdivide the interval in two and then use the quadrature on the half-intervals. Indeed, if the error is of the order  $h^k$ , the subdivision would lead to reduced error,  $2(h/2)^k < h^k$ , if  $k > 1$ .

An *adaptive quadrature* is an algorithm where the integration interval is subdivided into adaptively refined subintervals until the given accuracy goal is reached.

Adaptive algorithms are usually built on pairs of quadrature rules – a higher order rule,

$$Q = \sum_i w_i f(x_i), \quad (7.44)$$

where  $w_i$  are the weights of the higher order rule and  $Q$  is the higher order estimate of the integral, and a lower order rule,

$$q = \sum_i v_i f(x_i), \quad (7.45)$$

where  $v_i$  are the weights of the lower order rule and  $q$  is the lower order estimate of the integral. The difference between the higher order rule and the lower order rule gives an estimate of the error,

$$\delta Q = |Q - q|. \quad (7.46)$$

The integration result is accepted, if the error  $\delta Q$  is smaller than tolerance,

$$\delta Q < \delta + \epsilon |Q|, \quad (7.47)$$

where  $\delta$  is the absolute accuracy goal and  $\epsilon$  is the relative accuracy goal of the integration.

If the error estimate is larger than tolerance, the interval is subdivided into two half-intervals and the procedure applies recursively to subintervals with the same relative accuracy goal  $\epsilon$  and rescaled absolute accuracy goal  $\delta/\sqrt{2}$ .

The points  $x_i$  are usually chosen such that the two quadratures use the same points, and that the points can be reused in the subsequent recursive steps. The reuse of the function evaluations made at the previous step of adaptive integration is very important for the efficiency of the algorithm. The equally-spaced abscissas naturally provide for such a reuse.

As an example, Table 7.3 shows an implementation of the described algorithm using

$$x_i = \left\{ \frac{1}{6}, \frac{2}{6}, \frac{4}{6}, \frac{5}{6} \right\} \text{ (easily reusable points),} \quad (7.48)$$

$$w_i = \left\{ \frac{2}{6}, \frac{1}{6}, \frac{1}{6}, \frac{2}{6} \right\} \text{ (trapezium rule),} \quad (7.49)$$

$$v_i = \left\{ \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4} \right\} \text{ (rectangle rule).} \quad (7.50)$$

During recursion the function values at the points #2 and #3 are inherited from the previous step and need not be recalculated.

The points and weights are cited for the rescaled integration interval  $[0, 1]$ . The transformation of the points and weights to the original interval  $[a, b]$  is given as

$$\begin{aligned} x_i &\rightarrow a + (b - a)x_i, \\ w_i &\rightarrow (b - a)w_i. \end{aligned} \quad (7.51)$$

This implementation calculates directly the Riemann sums and can therefore deal with integrable singularities, although rather inefficiently.

More efficient adaptive routines keep track of the subdivisions of the interval and the local errors [14]. This allows detection of singularities and switching in their vicinity to specifically tuned quadratures. It also allows better estimates of local and global errors.

## 7.6 Variable transformation quadratures

The idea behind *variable transformation quadratures* is to apply the given quadrature—either with optimized or regularly spaced nodes—not to the original integral, but to

Table 7.3: Recursive adaptive integrator in C

```

#include<math.h>
#include<assert.h>
#include<stdio.h>
double adapt24(double f(double),double a, double b,
double acc, double eps, double f2, double f3, int nrec){
    assert(nrec < 1000000);
    double f1=f(a+(b-a)/6), f4=f(a+5*(b-a)/6);
    double Q=(2*f1+f2+f3+2*f4)/6*(b-a), q=(f1+f4+f2+f3)/4*(b-a);
    double tolerance=acc+eps*fabs(Q), error=fabs(Q-q);
    if(error < tolerance) return Q;
    else {
        double Q1=adapt24(f,a,(a+b)/2,acc/sqrt(2.),eps,f1,f2,nrec+1);
        double Q2=adapt24(f,(a+b)/2,b,acc/sqrt(2.),eps,f3,f4,nrec+1);
        return Q1+Q2; }
}
double adapt(
double f(double),double a,double b, double acc,double eps){
    double f2=f(a+2*(b-a)/6),f3=f(a+4*(b-a)/6); int nrec=0;
    return adapt24(f,a,b,acc,eps,f2,f3,nrec);
}
int main(){ //uses gcc nested functions
    int calls=0; double a=0,b=1,acc=0.001,eps=0.001;
    double f(double x){ calls++; return 1/sqrt(x);}; //nested function
    double Q=adapt(f,a,b,acc,eps); printf("Q=%g - calls=%d\n",Q,calls);
    return 0 ; }

```

a variable transformed integral [19],

$$\int_a^b f(x)dx = \int_{t_a}^{t_b} f(g(t))g'(t)dt \approx \sum_{i=1}^N w_i f(g(t_i))g'(t_i), \quad (7.52)$$

where the transformation  $x = g(t)$  is chosen such that the transformed integral better suits the given quadrature. Here  $g'$  denotes the derivative and  $[t_a, t_b]$  is the corresponding interval in the new variable.

For example, the Gauss-Legendre quadrature assumes the integrand can be well represented with polynomials and performs poorly on integrals with integrable singularities like

$$I = \int_0^1 \frac{1}{2\sqrt{x}} dx. \quad (7.53)$$

However, a simple variable transformation  $x = t^2$  removes the singularity,

$$I = \int_0^1 dt, \quad (7.54)$$

and the Gauss-Legendre quadrature for the transformed integral gives exact result. The price is that the transformed quadrature performs less effectively on smooth functions.

Some of the popular variable transformation quadratures are Clenshaw-Curtis [8], based on the transformation

$$\int_{-1}^1 f(x)dx = \int_0^\pi f(\cos \theta) \sin \theta d\theta, \quad (7.55)$$

and “tanh-sinh” quadrature [19], based on the transformation

$$\int_{-1}^1 f(x)dx = \int_{-\infty}^{\infty} f\left(\tanh\left(\frac{\pi}{2}\sinh(t)\right)\right) \frac{\pi}{2} \frac{\cosh(t)}{\cosh^2\left(\frac{\pi}{2}\sinh(t)\right)} dt. \quad (7.56)$$

Generally, the equally spaced trapezium rule is used after the transformation.

## 7.7 Infinite intervals

One way to calculate an integral over infinite interval is to transform it by a variable substitution into an integral over a finite interval. The latter can then be evaluated by ordinary integration methods. Table 7.4 lists several of such transformation.

Table 7.4: Variable transformations reducing infinite interval integrals into integrals over finite intervals.

$$\int_{-\infty}^{+\infty} f(x)dx = \int_{-1}^{+1} f\left(\frac{t}{1-t^2}\right) \frac{1+t^2}{(1-t^2)^2} dt, \quad (7.57)$$

$$\int_{-\infty}^{+\infty} f(x)dx = \int_0^1 \left( f\left(\frac{1-t}{t}\right) + f\left(-\frac{1-t}{t}\right) \right) \frac{dt}{t^2}, \quad (7.58)$$

$$\int_a^{+\infty} f(x)dx = \int_0^1 f\left(a + \frac{t}{1-t}\right) \frac{1}{(1-t)^2} dt, \quad (7.59)$$

$$\int_a^{+\infty} f(x)dx = \int_0^1 f\left(a + \frac{1-t}{t}\right) \frac{dt}{t^2}, \quad (7.60)$$

$$\int_{-\infty}^b f(x)dx = \int_{-1}^0 f\left(b + \frac{t}{1+t}\right) \frac{1}{(1+t)^2} dt, \quad (7.61)$$

$$\int_{-\infty}^b f(x)dx = \int_0^1 f\left(b - \frac{1-t}{t}\right) \frac{dt}{t^2}. \quad (7.62)$$


---



# Chapter 8

## Minimization and optimization

### 8.1 Introduction

*Minimization (maximization)* is a problem of finding the minimum (maximum) of a given — generally non-linear — real valued function  $\phi(\mathbf{x})$  (often called the *objective function*) of an  $n$ -dimensional argument  $\mathbf{x} \doteq \{x_1, \dots, x_n\}$ .

Minimization is a simple case of a more general problem — *optimization* — which includes finding best available values of the objective function within a given domain and subject to given constraints.

Minimization is not unrelated to root-finding: at the minimum all partial derivatives of the objective function vanish,

$$\frac{\partial \phi}{\partial x_i} = 0 \Big|_{i=1, \dots, n}, \quad (8.1)$$

and one can alternatively solve this system of (non-linear) equations.

### 8.2 Local minimization

#### 8.2.1 Newton's methods

Newton's method is based on the quadratic approximation of the objective function  $\phi$  in the vicinity of the suspected minimum,

$$\phi(\mathbf{x} + \Delta \mathbf{x}) \approx \phi(\mathbf{x}) + \nabla \phi(\mathbf{x})^\top \Delta \mathbf{x} + \frac{1}{2} \Delta \mathbf{x}^\top \mathbf{H}(\mathbf{x}) \Delta \mathbf{x}, \quad (8.2)$$

where the vector  $\nabla\phi(\mathbf{x})$  is the gradient of the objective function at the point  $\mathbf{x}$ ,

$$\nabla\phi(\mathbf{x}) \doteq \left\{ \frac{\partial\phi(\mathbf{x})}{\partial x_i} \right\}_{i=1,\dots,n}, \quad (8.3)$$

and  $\mathbf{H}(\mathbf{x})$  is the *Hessian matrix* – a square matrix of second-order partial derivatives of the objective function at the point  $\mathbf{x}$ ,

$$\mathbf{H}(\mathbf{x}) \doteq \left\{ \frac{\partial^2\phi(\mathbf{x})}{\partial x_i\partial x_j} \right\}_{i,j\in 1,\dots,n}. \quad (8.4)$$

The minimum of the quadratic form (8.2), as function of  $\Delta\mathbf{x}$ , is found at the point where its gradient with respect to  $\Delta\mathbf{x}$  vanishes,

$$\nabla\phi(\mathbf{x}) + \mathbf{H}(\mathbf{x})\Delta\mathbf{x} = 0. \quad (8.5)$$

This gives an approximate step towards the minimum, called the *Newton's step*,

$$\Delta\mathbf{x} = -\mathbf{H}(\mathbf{x})^{-1}\nabla\phi(\mathbf{x}). \quad (8.6)$$

The original Newton's method is simply the iteration,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}(\mathbf{x}_k)^{-1}\nabla\phi(\mathbf{x}_k), \quad (8.7)$$

where at each iteration the full Newton's step is taken and the Hessian matrix is recalculated.

In practice, instead of calculating  $\mathbf{H}^{-1}$  one rather solves the linear equation (8.5).

Usually the Newton's method is modified to take a smaller step  $\mathbf{s}$ ,

$$\mathbf{s} = \lambda\Delta\mathbf{x}, \quad (8.8)$$

with  $0 < \lambda < 1$ . The factor  $\lambda$  can be found by a backtracking algorithm similar to that in the Newton's method for root-finding. One starts with  $\lambda = 1$  and then backtracks,  $\lambda \leftarrow \lambda/2$ , until the *Armijo condition*,

$$\phi(\mathbf{x} + \mathbf{s}) < \phi(\mathbf{x}) + \alpha\mathbf{s}^T\nabla\phi(\mathbf{x}), \quad (8.9)$$

is satisfied (or the minimal  $\lambda$  is reached, in which case the step is taken unconditionally). The parameter  $\alpha$  can be chosen as small as  $10^{-4}$ .

## 8.2.2 Quasi-Newton methods

Quasi-Newton methods are variations of the Newton's method which attempt to avoid recalculation of the Hessian matrix at each iteration, trying instead certain updates

based on the analysis of the gradient vectors. The update  $\delta\mathbf{H}$  is usually chosen to satisfy the condition

$$\nabla\phi(\mathbf{x} + \mathbf{s}) = \nabla\phi(\mathbf{x}) + (\mathbf{H} + \delta\mathbf{H})\mathbf{s}, \quad (8.10)$$

called *secant equation*, which is the Taylor expansion of the gradient.

The secant equation is under-determined in more than one dimension as it consists of only  $n$  equations for the  $n^2$  unknown elements of the update  $\delta\mathbf{H}$ . Various quasi-Newton methods use different choices for the form of the solution of the secant equation.

In practice one typically uses the inverse Hessian matrix (often—but not always—denoted as  $\mathbf{B}$ ) and applies the updates directly to the inverse matrix thus avoiding the need to solve the linear equation (8.5) at each iteration.

For the inverse Hessian matrix the secant equation (8.10) reads

$$(\mathbf{B} + \delta\mathbf{B})\mathbf{y} = \mathbf{s}, \quad (8.11)$$

or, in short,

$$\delta\mathbf{B}\mathbf{y} = \mathbf{u}, \quad (8.12)$$

where  $\mathbf{B} \doteq \mathbf{H}^{-1}$ ,  $\mathbf{y} \doteq \nabla\phi(\mathbf{x} + \mathbf{s}) - \nabla\phi(\mathbf{x})$ , and  $\mathbf{u} \doteq \mathbf{s} - \mathbf{B}\mathbf{y}$ .

One usually starts with the identity matrix as the zeroth approximation for the inverse Hessian matrix and then applies the updates.

If the minimal  $\lambda$  is reached during the backtracking line-search—which might be a signal of lost precision in the approximate (inverse) Hessian matrix—it is advisable to reset the current inverse Hessian matrix to identity matrix.

### Broyden's update

The Broyden's update is chosen in the form

$$\delta\mathbf{B} = \mathbf{c}\mathbf{s}^T. \quad (8.13)$$

where the vector  $\mathbf{c}$  is found from the condition (8.12),

$$\mathbf{c} = \frac{\mathbf{u}}{\mathbf{s}^T\mathbf{y}}. \quad (8.14)$$

Sometimes the dot-product  $\mathbf{s}^T\mathbf{y}$  becomes very small or even zero which results in serious numerical difficulties. One can avoid this by only performing update if the condition  $|\mathbf{s}^T\mathbf{y}| > \epsilon$  is satisfied where  $\epsilon$  is a small number, say  $10^{-6}$ .

### Symmetric Broyden's update

The Broyden's update (8.13) is not symmetric (while the Hessian matrix should be) which is an obvious drawback. Therefore a better approximation might be the symmetric Broyden's update,

$$\delta\mathbf{B} = \mathbf{a}\mathbf{s}^\top + \mathbf{s}\mathbf{a}^\top. \quad (8.15)$$

The vector  $\mathbf{a}$  is again found from the condition (8.12),

$$\mathbf{a} = \frac{\mathbf{u} - \gamma\mathbf{s}}{\mathbf{s}^\top\mathbf{y}}, \quad (8.16)$$

where  $\gamma = (\mathbf{u}^\top\mathbf{y})/(2\mathbf{s}^\top\mathbf{y})$ .

Again one only performs the update if  $|\mathbf{s}^\top\mathbf{y}| > \epsilon$ .

### SR1 update

The symmetric-rank-1 update (SR1) is chosen in the form

$$\delta\mathbf{B} = \mathbf{v}\mathbf{v}^\top, \quad (8.17)$$

where the vector  $\mathbf{v}$  is again found from the condition (8.10), which gives

$$\delta\mathbf{B} = \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{y}}. \quad (8.18)$$

Again, one only performs the update if denominator is not too small, that is,  $|\mathbf{u}^\top\mathbf{y}| > \epsilon$ .

### Other popular updates

The wikipedia article “Quasi-Newton method” list several other popular updates.

### 8.2.3 Downhill simplex method

The *downhill simplex method* [17] (also called “Nelder-Mead” or “amoeba”) is a commonly used minimization algorithm where the minimum of a function in an  $n$ -dimensional space is found by transforming a simplex—a polytope with  $n+1$  vertexes—according to the function values at the vertexes, moving it downhill until it converges towards the minimum.

The advantages of the downhill simplex method is its stability and the lack of use of derivatives. However, the convergence is relatively slow as compared to Newton's methods.

In order to introduce the algorithm we need the following definitions:

- Simplex: a figure (polytope) represented by  $n+1$  points, called vertexes,  $\{\mathbf{p}_1, \dots, \mathbf{p}_{n+1}\}$  (where each point  $\mathbf{p}_k$  is an  $n$ -dimensional vector).
- Highest point: the vertex,  $\mathbf{p}_{\text{hi}}$ , with the highest value of the function:  $\phi(\mathbf{p}_{\text{hi}}) = \max_k \phi(\mathbf{p}_k)$ .
- Lowest point: the vertex,  $\mathbf{p}_{\text{lo}}$ , with the lowest value of the function:  $\phi(\mathbf{p}_{\text{lo}}) = \min_k \phi(\mathbf{p}_k)$ .
- Centroid: the center of gravity of all points, except for the highest:  $\mathbf{p}_{\text{ce}} = \frac{1}{n} \sum_{(k \neq \text{hi})} \mathbf{p}_k$

The simplex is moved downhill by a combination of the following elementary operations:

1. Reflection: the highest point is reflected against the centroid,  $\mathbf{p}_{\text{hi}} \rightarrow \mathbf{p}_{\text{re}} = \mathbf{p}_{\text{ce}} + (\mathbf{p}_{\text{ce}} - \mathbf{p}_{\text{hi}})$ .
2. Expansion: the highest point reflects and then doubles its distance from the centroid,  $\mathbf{p}_{\text{hi}} \rightarrow \mathbf{p}_{\text{ex}} = \mathbf{p}_{\text{ce}} + 2(\mathbf{p}_{\text{ce}} - \mathbf{p}_{\text{hi}})$ .
3. Contraction: the highest point halves its distance from the centroid,  $\mathbf{p}_{\text{hi}} \rightarrow \mathbf{p}_{\text{co}} = \mathbf{p}_{\text{ce}} + \frac{1}{2}(\mathbf{p}_{\text{hi}} - \mathbf{p}_{\text{ce}})$ .
4. Reduction: all points, except for the lowest, move towards the lowest points halving the distance.  $\mathbf{p}_{k \neq \text{lo}} \rightarrow \frac{1}{2}(\mathbf{p}_k + \mathbf{p}_{\text{lo}})$ .

Table 8.1 shows one possible algorithm for the downhill simplex method, and a C-implementation of simplex operations and the amoeba algorithm can be found in Table 8.4 and Table 8.4.

## 8.3 Global optimization

Global optimization is the problem of locating (a good approximation to) the global minimum of a given objective function in a search space large enough to prohibit exhaustive enumeration.

When only a small sub-space of the search space can be realistically sampled the stochastic methods usually come to the fore.

A good local minimizer converges to the nearest local minimum relatively fast, so one possible global minimizer can be constructed by simply starting the local minimizer from different random starting points.

In the following several popular global minimization algorithms are shortly described.

Table 8.1: Downhill simplex (Nelder-Mead) algorithm

```

REPEAT :
  find highest, lowest, and centroid points of the simplex
  try reflection
  IF  $\phi(\text{reflected}) < \phi(\text{lowest})$  :
    try expansion
    IF  $\phi(\text{expanded}) < \phi(\text{reflected})$  :
      accept expansion
    ELSE :
      accept reflection
  ELSE :
    IF  $\phi(\text{reflected}) < \phi(\text{highest})$  :
      accept reflection
    ELSE :
      try contraction
      IF  $\phi(\text{contracted}) < \phi(\text{highest})$  :
        accept contraction
      ELSE :
        do reduction
UNTIL converged (e.g. size(simplex) < tolerance)

```

### 8.3.1 Simulated annealing

Simulated annealing is a stochastic metaheuristic algorithm for global minimization. The name and inspiration come from annealing—heating up and cooling slowly—in material science. The slow cooling allows a piece of material to reach a state with “lowest energy”.

The objective function in the space of states is interpreted as some sort of potential energy and the states—the points in the search space—are interpreted as physical states of some physical system. The system attempts to make transitions from its current state to some randomly sampled nearest states with the goal to eventually reach the state with minimal energy – the global minimum.

The system is attached to a thermal reservoir with certain temperature. Each time the energy of the system is measured the reservoir supplies it with a random amount of thermal energy sampled from the Boltzmann distribution,

$$P(E) = Te^{-E/T} . \quad (8.19)$$

If the temperature equals zero the system can only make transitions to the neighboring states with lower potential energy. In this case the algorithm turns merely into a local minimizer with random sampling.

If temperature is finite the system is able to climb up the ridges of the potential

Table 8.2: Simulated annealing algorithm

```

state ← start_state
T ← start_temperature
energy ← E(state)
REPEAT :
  new_state ← neighbour(state)
  new_energy ← E(new_state)
  IF new_energy < energy :
    state ← new_state
    energy ← new_energy
  ELSE :
    do with probability  $\exp\left(-\frac{\text{new\_energy}-\text{energy}}{T}\right)$  :
      state ← new_state
      energy ← new_energy
  reduce_temperature_according_to_schedule(T)
UNTIL terminated

```

energy—about as high as the current temperature—and thus escape from local minima and hopefully eventually reach the global minimum.

One typically starts the simulation with some finite temperature on the order of the height of the typical hills of the potential energy surface, letting the system to wander almost unhindered around the landscape with a good chance to locate if not the best then at least a good enough minimum. The temperature is then slowly reduced following some annealing schedule which may be supplied by the user but must end with  $T = 0$  towards the end of the allotted time budget.

Table 8.2 lists one possible variant of the algorithm.

The function `neighbour(state)` should return a randomly chosen neighbour of the given state.

Downhill simplex method can incorporate simulated annealing by adding the stochastic thermal energy to the values of the objective function at the vertices.

### 8.3.2 Quantum annealing

Quantum annealing is a general global minimization algorithm which—like simulated annealing—also allows the search path to escape from local minima. However instead of the thermal jumps over the potential barriers quantum annealing allows the system to tunnel through the barriers.

In its simple incarnation the quantum annealing algorithm allows the system to attempt transitions not only to the nearest states but also to distant states within

Table 8.3: Quantum annealing algorithm

```

state ← start_state
energy ← E(state)
R ← start_radius
REPEAT :
  new_state ← random_neighbour_within_radius(state ,R)
  new_energy ← E(new_state)
  IF new_energy < energy :
    state ← new_state
    energy ← new_energy
  reduce_radius_according_to_schedule(R)
UNTIL terminated

```

certain "tunneling distance" from the current state. The transition is accepted only if it reduces the potential energy of the system.

At the beginning of the minimization procedure the tunnelling distance is large—on the order of the size of the region where the global minimum is suspected to be located—allowing the system to explore the region. The tunneling distance is then slowly reduced according to a schedule such that by the end of the allotted time the tunnelling distance reduces to zero at which point the system hopefully is in the state with minimal energy.

### 8.3.3 Evolutionary algorithms

Unlike annealing algorithms, which follow the motion of only one point in the search space, the evolutionary algorithms typically follow a set of points called a population of individuals. A bit like the downhill simplex method which follows the motion of a set of points – the simplex.

The population evolves towards more fit individuals where fitness is understood in the sense of minimizing the objective functions. The parameters of the individuals (for example, the coordinates of the points in multi-dimensional minimization of a continuous objective function) are called genes.

The algorithm proceeds iteratively with the population in each iteration called a generation. In each generation the fitness of each individual—typically, the value of the objective function—is evaluated and the new generation is generated stochastically from the gene pool of the current generation through crossovers and mutations such that the genes of more fit individuals have a better chance of propagating into the next generation.

Each new individual in the next generation is produced from a pair of "parent"



individuals of the current generation. The use of two "parents" is biologically inspired, in practice more than two "parents" can be used as well. The parents for a new individual are selected from the individuals of the current generation through a fitness based stochastic process where fitter individuals are more likely to be selected.

The "child" individual shares many characteristics of its "parents". In the simplest case the "child" may get its genes by simply averaging the genes of its parents. Then a certain amount of mutations—random changes in the genes—are added to the "child's" genes.

Generation of "children" continues until the population of the new generation reaches the appropriate size after which the iteration repeats itself.

The algorithm is terminated when the fitness level of the population is deemed sufficient or when the allocated budget is exhausted.

## 8.4 Implementation in C

Table 8.4: C implementation of simplex operations

```

void reflection
(double* highest, double* centroid, int dim, double* reflected){
    for(int i=0;i<dim;i++) reflected[i]=2*centroid[i]-highest[i];
}
void expansion
(double* highest, double* centroid, int dim, double* expanded) {
    for(int i=0;i<dim;i++) expanded[i]=3*centroid[i]-2*highest[i];
}
void contraction
(double* highest, double* centroid, int dim, double* contracted){
    for(int i=0;i<dim;i++)
        contracted[i]=0.5*centroid[i]+0.5*highest[i];
}
void reduction(double** simplex, int dim, int lo){
    for(int k=0;k<dim+1;k++) if(k!=lo) for(int i=0;i<dim;i++)
        simplex[k][i]=0.5*(simplex[k][i]+simplex[lo][i]);
}
double distance(double* a, double* b, int dim){
    double s=0; for(int i=0;i<dim;i++) s+=pow(b[i]-a[i],2);
    return sqrt(s);
}
double size(double** simplex, int dim){
    double s=0; for(int k=1;k<dim+1;k++){
        double dist=distance(simplex[0], simplex[k], dim);
        if(dist>s) s=dist; }
    return s;
}

```

```

void simplex_update(double** simplex, double* f_values, int d,
int* hi, int* lo, double* centroid) {
    *hi=0; *lo=0; double highest=f_values[0], lowest =f_values[0];
    for(int k=1;k<d+1;k++) {
        double next=f_values[k];
        if(next>highest){ highest=next; *hi=k;}
        if(next<lowest) { lowest=next; *lo=k;} }
    for(int i=0;i<d;i++) {
        double s=0; for(int k=0;k<d+1;k++) if(k!=*hi) s+=simplex[k][i];
        centroid[i]=s/d; }
}
void simplex_initiate(
double fun(double*), double** simplex, double* f_values, int d,
int* hi, int* lo, double* centroid) {
    for(int k=0;k<d+1;k++) f_values[k]=fun(simplex[k]);
    simplex_update(simplex, f_values, d, hi, lo, centroid);
}

```

Table 8.5: C implementation of downhill simplex algorithm

```

int downhill_simplex(
    double F(double*),double**simplex,int d,double simplex_size_goal)
{
int hi,lo,k=0; double centroid[d], F_value[d+1], p1[d], p2[d];
simplex_initiate(F,simplex,F_value,d,&hi,&lo,centroid);
while(size(simplex,d)>simplex_size_goal){
    simplex_update(simplex,F_value,d,&hi,&lo,centroid);
    reflection(simplex[hi],centroid,d,p1); double f_re=F(p1);
    if(f_re<F_value[lo]){ // reflection looks good: try expansion
        expansion(simplex[hi],centroid,d,p2); double f_ex=F(p2);
        if(f_ex<f_re){ // accept expansion
            for(int i=0;i<d;i++)simplex[hi][i]=p2[i]; F_value[hi]=f_ex;}
        else{ // reject expansion and accept reflection
            for(int i=0;i<d;i++)simplex[hi][i]=p1[i]; F_value[hi]=f_re;}}
    else{ // reflection wasn't good
        if(f_re<F_value[hi]){ // ok, accept reflection
            for(int i=0;i<d;i++)simplex[hi][i]=p1[i]; F_value[hi]=f_re;}
        else{ // try contraction
            contraction(simplex[hi],centroid,d,p1); double f_co=F(p1);
            if(f_co<F_value[hi]){ // accept contraction
                for(int i=0;i<d;i++)simplex[hi][i]=p1[i]; F_value[hi]=f_co;}
            else{ // do reduction
                reduction(simplex,d,lo);
                simplex_initiate(F,simplex,F_value,d,&hi,&lo,centroid);}}}
    k++;} return k;
}

```



## Chapter 9

# Monte Carlo integration

### 9.1 Introduction

*Monte Carlo integration* is a quadrature (cubature) where the nodes are chosen randomly. Typically no assumption is made about the smoothness of the integrand, not even that it is continuous.

Monte Carlo algorithms are particularly suited for multi-dimensional integrations where one of the problems is that the integration region,  $\Omega$ , might have a quite complicated boundary which can not be easily described by simple functions. On the other hand, it is usually much easier to find out whether a given point lies within the integration region or not. Therefore a popular strategy is to create an auxiliary rectangular volume,  $V$ , which encompasses the integration volume  $\Omega$ , and an auxiliary function which coincides with the integrand inside the volume  $\Omega$  and is equal zero outside. Then the integral of the auxiliary function over the auxiliary volume is equal the original integral.

However, the auxiliary function is generally non-continuous at the boundary; thus ordinary quadratures—that assume continuity of the integrand—are bound to have difficulties here. On the contrary the Monte-Carlo quadratures will do just as good (or as bad) as with continuous integrands.

A typical implementation of a Monte Carlo algorithm integrates the given function over a rectangular volume, specified by the coordinates of its "lower-left" and "upper-right" vertices, assuming the user has provided the encompassing volume with the auxiliary function.

Plain Monte Carlo algorithm distributes points uniformly throughout the integration region using uncorrelated pseudo-random sequences of points.

Adaptive algorithms, such as VEGAS and MISER, distribute points non-uniformly

in an attempt to reduce integration error using correspondingly *importance* and *stratified* sampling.

Yet another strategy to reduce the error is to use correlated quasi-random sequences.

The GNU Scientific Library, GSL, implements a plain Monte Carlo integration algorithm; a stratified sampling algorithm, MISER; an importance sampling algorithm, VEGAS; and a number of quasi-random generators.

## 9.2 Plain Monte Carlo sampling

Plain Monte Carlo is a quadrature with random abscissas and equal weights,

$$\int_V f(\mathbf{x})dV \approx w \sum_{i=1}^N f(\mathbf{x}_i) , \quad (9.1)$$

where  $\mathbf{x}$  is a point in the multi-dimensional integration space. One free parameter,  $w$ , allows one condition to be satisfied: the quadrature must integrate exactly a constant function. This gives  $w = V/N$ ,

$$\int_V f(\mathbf{x})dV \approx \frac{V}{N} \sum_{i=1}^N f(\mathbf{x}_i) \doteq V \langle f \rangle . \quad (9.2)$$

Under the assumptions of the *central limit theorem* the error of the integration can be estimated as

$$\text{error} = V \frac{\sigma}{\sqrt{N}} , \quad (9.3)$$

where  $\sigma$  is the variance of the sample,

$$\sigma^2 = \langle f^2 \rangle - \langle f \rangle^2 . \quad (9.4)$$

The familiar  $1/\sqrt{N}$  convergence of a random walk process is quite slow: to reduce the error by a factor 10 requires 100-fold increase in the number of sample points.

Expression (9.3) provides only a statistical estimate of the error, which is not a strict bound; random sampling may not uncover all the important features of the function, resulting in an underestimate of the error.

A simple implementation of the plain Monte Carlo algorithm is shown in Table 9.1.

Table 9.1: Plain Monte Carlo integrator

```

#include <math.h>
#include <stdlib.h>
#define RND ((double)rand()/RAND_MAX)
void randomx(int dim, double *a, double *b, double *x)
{ for(int i=0;i<dim;i++) x[i]=a[i]+RND*(b[i]-a[i]); }

void plainmc(int dim, double *a, double *b,
double f(double* x),int N, double*result , double*error)
{ double V=1; for(int i=0;i<dim;i++) V*=b[i]-a[i];
double sum=0, sum2=0, fx , x[dim];
for(int i=0;i<N;i++){ randomx(dim,a,b,x); fx=f(x);
sum+=fx; sum2+=fx*fx; }
double avr = sum/N, var = sum2/N-avr*avr;
*result = avr*V; *error = sqrt(var/N)*V;
}

```

### 9.3 Importance sampling

Suppose the points are distributed not uniformly but with some density  $\rho(\mathbf{x})$ . That is, the number of points  $\Delta n$  in the volume  $\Delta V$  around point  $\mathbf{x}$  is given as

$$\Delta n = \frac{N}{V} \rho(\mathbf{x}) \Delta V, \quad (9.5)$$

where  $\rho$  is normalised such that  $\int_V \rho dV = V$ .

The estimate of the integral is then given as

$$\int_V f(\mathbf{x}) dV \approx \sum_{i=1}^N f(\mathbf{x}_i) \Delta V_i = \sum_{i=1}^N f(\mathbf{x}_i) \frac{V}{N \rho(\mathbf{x}_i)} = V \left\langle \frac{f}{\rho} \right\rangle, \quad (9.6)$$

where

$$\Delta V_i = \frac{V}{N \rho(x_i)} \quad (9.7)$$

is the volume-per-point at the point  $x_i$ .

The corresponding variance is now given by

$$\sigma^2 = \left\langle \left( \frac{f}{\rho} \right)^2 \right\rangle - \left\langle \frac{f}{\rho} \right\rangle^2. \quad (9.8)$$

Apparently if the ratio  $f/\rho$  is close to a constant, the variance is reduced.

It is tempting to take  $\rho = |f|$  and sample directly from the integrand. However in practice evaluations of the integrand are typically expensive. Therefore a better strategy is to build an approximate density in the product form,  $\rho(x, y, \dots, z) = \rho_x(x)\rho_y(y)\dots\rho_z(z)$ , and then sample from this approximate density. A popular routine of this sort is called VEGAS.

## 9.4 Stratified sampling

Stratified sampling is a generalisation of the recursive adaptive integration algorithm to random quadratures in multi-dimensional spaces.

Table 9.2: Recursive stratified sampling algorithm

```

sample  $N$  random points with plain Monte Carlo;
estimate the average and the error;
IF the error is acceptable :
    RETURN the average and the error;
ELSE :
    FOR EACH dimension :
        subdivide the volume in two along the dimension;
        estimate the sub-variances in the two sub-volumes;
    pick the dimension with the largest sub-variance;
    subdivide the volume in two along this dimension;
    dispatch two recursive calls to each of the sub-volumes;
    estimate the grand average and grand error;
    RETURN the grand average and grand error;

```

The ordinary “dividing by two” strategy does not work for multi-dimensional integrations as the number of sub-volumes grows way too fast to keep track of. Instead one estimates along which dimension a subdivision should bring the most dividends and only subdivides along this dimension. Such strategy is called *recursive stratified sampling*. A simple variant of this algorithm is presented in Table 9.2.

In a stratified sample the points are concentrated in the regions where the variance of the function is largest, see the illustration in Figure 9.1.

The naive implementation in Table 9.2 keeps throwing points until the given tolerance is achieved. This might be problematic for multidimensional integrals where for the given tolerance one might need excessively many points. The GSL’s implementation instead takes the number of points as the argument and then returns the estimate of the integral and the estimate of the error.



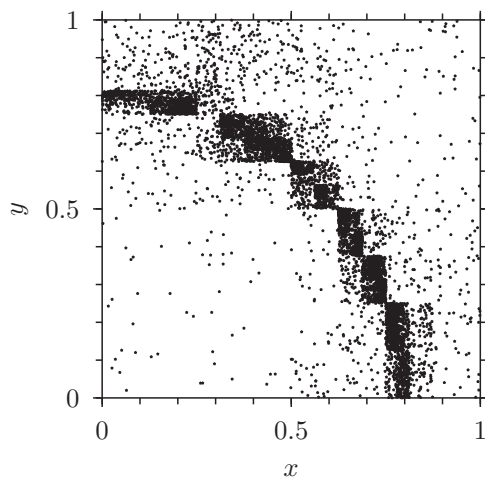


Figure 9.1: Stratified sample of a discontinuous function,  $f(x, y) = 1$  if  $x^2 + y^2 < 0.8^2$  otherwise  $f(x, y) = 0$ , built with the algorithm in Table 9.2.

## 9.5 Quasi-random (low-discrepancy) sampling

Pseudo-random sampling has high discrepancy<sup>1</sup>: it typically creates regions with high density of points and other regions with low density of points, see an illustration on Figure 9.2 (left). With pseudo-random sampling there is a finite probability that all the  $N$  points would fall into one half of the region and none into the other half.

*Quasi-random* sequences avoid this phenomenon by distributing points in a highly correlated manner with a specific requirement of low discrepancy, see Figure 9.2 for an example. Quasi-random sampling is like a computation on a grid where the grid constant must not be known in advance as the grid is ever gradually refined and the points are always distributed uniformly over the region. The computation can be stopped at any time.

By placing points more evenly than at random, the quasi-random sequences try to improve on the  $1/\sqrt{N}$  convergence rate of pseudo-random sampling.

The central limit theorem does not apply in this case as the points are not statistically independent. Therefore the variance can not be used as an estimate of the error. The error estimation is actually not trivial. In practice one can employ two different sequences and use the difference in the resulting integrals as an error estimate.

---

<sup>1</sup>discrepancy is a measure of how unevenly the points are distributed over the region.

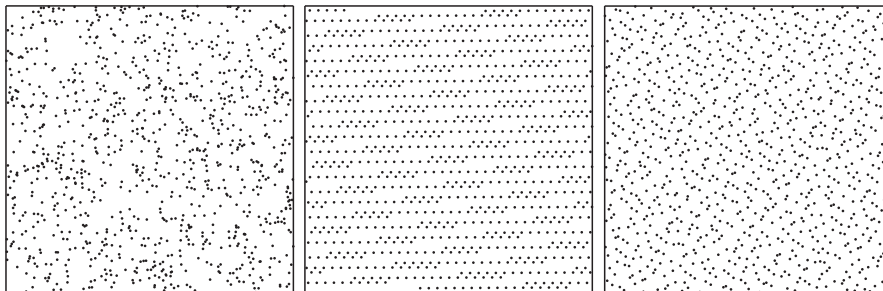


Figure 9.2: Typical distributions of pseudo-random points (left), and quasi-random low-discrepancy points: lattice (center) and base-2/3 Halton (right) sequences. The first thousand points are plotted in each case.

### 9.5.1 Van der Corput and Halton sequences

A *van der Corput sequence* is a low-discrepancy sequence over the unit interval. It is constructed by reversing the base- $b$  representation of the sequence of natural numbers  $(1, 2, 3, \dots)$ . For example, the decimal van der Corput sequence begins as

$$0.1, 0.2, 0.3, \dots, 0.8, 0.9, 0.01, 0.11, 0.21, 0.31, \dots, 0.91, 0.02, 0.12, \dots \quad (9.9)$$

In a base- $b$  representation a natural number  $n$  with  $s$  digits  $\{d_i \mid i = 1 \dots s, 0 \leq d_i < b\}$  is given as

$$n = \sum_{k=1}^s d_k b^{k-1}. \quad (9.10)$$

The corresponding base- $b$  van der Corput number  $q_b(n)$  is then given as

$$q_b(n) = \sum_{k=1}^s d_k b^{-k}. \quad (9.11)$$

Here is a C implementation of this algorithm,

```
double corput(int n, int base){
  double q=0, bk=(double)1/base;
  while(n>0){ q += (n % base)*bk; n /= base; bk /= base; }
  return q; }
```

The van der Corput numbers of any base are uniformly distributed over the unit interval. They also form a dense set in the unit interval: there exists a subsequence of the van der Corput sequence which converges to any given real number in  $[0, 1]$ .

The Halton sequence is a generalization of the van der Corput sequence to  $d$ -dimensional spaces. One chooses a set of coprime bases  $b_1, \dots, b_d$  and then for each dimension  $i$  generates a van der Corput sequence with its own base  $b_i$ . The  $n$ -th Halton  $d$ -dimensional point  $\mathbf{x}$  in the unit volume is then given as

$$\mathbf{x}_{b_1, \dots, b_d}(n) = \{q_{b_1}(n), \dots, q_{b_d}(n)\}. \quad (9.12)$$

Here is a C implementation which calls the `corput` function listed above,

```
#include<assert.h>
void halton(int n, int d, double *x){
    int base[]={2,3,5,7,11,13,17,19,23,29,31,37,41,43,47,53,59,61,67};
    int maxd=sizeof(base)/sizeof(int); assert(d <= maxd);
    for(int i=0;i<d;i++) x[i]=corput(n, base[i]); }
```

### 9.5.2 Lattice rules

In the simplest incarnation a lattice rule can be defined as follows. Let the *generating vector*  $\mathbf{z} = \{\alpha_i \mid i = 1, \dots, d\}$  — where  $d$  is the dimension of the integration space — be a set of cleverly chosen irrational numbers. Then the  $n$ -th point (in the unit volume) of the sequence is given as

$$\mathbf{x}(n) = \text{frac}(n\mathbf{z}) \equiv \{\text{frac}(n\alpha_1), \dots, \text{frac}(n\alpha_d)\}, \quad (9.13)$$

where  $\text{frac}(x)$  is the fractional part of  $x$ .

An implementation of this algorithm in C is given in Table 9.3 and an illustration of such sequence is shown on Figure 9.2 (center).

## 9.6 Implementations

Table 9.3: Lattice low-discrepancy quasi-random sequence in C.

```

#define fracl(x) ((x)-floorl(x))
#define PI 3.1415926535897932384626433832795028841971693993751L
#define real long double
void lattice(int d, double *x){
    static int dim=0, n=-1; static real *alpha; int i;
    if(d<0){ /* d<0 is the signal to (re-)initialize the lattice */
        dim=d; n=0; alpha=(real*)realloc(alpha,dim*sizeof(real));
        for(i=0;i<dim;i++) alpha[i]=fracl(sqrtl(PI+i));
    }
    else if(d>0){
        n++; assert(d==dim && n>0);
        for(i=0;i<dim;i++)x[i]=fracl(n*alpha[i]);
    }
    else if(alpha!=NULL) free(alpha);
    return;
}

```

Table 9.4: C implementation of the stratified sampling algorithm

```

#include<math.h>
#include<stdlib.h>
#include<stdio.h>
#define FOR(k) for (int k=0;k<dim;k++)
#define RND (double)rand()/RAND_MAX

double strata (
    int dim, double f(int dim, double*x),
    double*a, double*b,
    double acc, double eps,
    int n_reuse, double mean_reuse)
{
    int N=16*dim;
    double V=1; FOR(k) V*=b[k]-a[k];
    int n_left [dim], n_right [dim];
    double x[dim], mean_left [dim], mean_right [dim], mean=0;
    FOR(k){ mean_left [k]=0; mean_right [k]=0; n_left [k]=0; n_right [k]=0; }
    for (int i=0; i<N; i++){
        FOR(k) x [k]=a [k]+RND*(b[k]-a [k]);
        double fx=f(dim, x);
        mean+=fx;
        FOR(k){
            if (x[k]>(a[k]+b[k])/2){ n_right [k]++; mean_right [k]+=fx; }
            else { n_left [k]++; mean_left [k]+=fx; }
        }
    }
    mean/=N;
    FOR(k){ mean_left [k]/=n_left [k]; mean_right [k]/=n_right [k]; }

    int kdiv=0; double maxvar=0;
    FOR(k){
        double var=fabs (mean_right [k]-mean_left [k]);
        if (var>maxvar){ maxvar=var; kdiv=k; }
    }

    double integ=(mean*N+mean_reuse*n_reuse)/(N+n_reuse)*V;
    double error=fabs (mean_reuse-mean)*V;
    double toler=acc+fabs (integ)*eps;
    if (error<toler) return integ;

    double a2 [dim], b2 [dim]; FOR(k) a2 [k]=a [k]; FOR(k) b2 [k]=b [k];
    a2 [kdiv]=(a [kdiv]+b [kdiv])/2; b2 [kdiv]=(a [kdiv]+b [kdiv])/2;
    double integ_left=
        strata (dim, f, a, b2, acc/sqrt (2), eps, n_left [kdiv], mean_left [kdiv]);
    double integ_right=
        strata (dim, f, a2, b, acc/sqrt (2), eps, n_right [kdiv], mean_right [kdiv]);
    return integ_left+integ_right;
}

```



## Chapter 10

# Power iteration methods and Krylov subspaces

### 10.1 Power iteration

Power method is an iterative method to calculate an eigenvalue and the corresponding eigenvector of a (real symmetric) matrix  $A$  using the *power iteration*

$$\mathbf{x}_{i+1} = \mathbf{A}\mathbf{x}_i . \quad (10.1)$$

The iteration converges to the eigenvector with the largest eigenvalue. Indeed, according to the spectral theorem the eigenvectors  $\mathbf{v}_i$ ,

$$\mathbf{A}\mathbf{v}_i = \lambda_i \mathbf{v}_i , \quad (10.2)$$

of a real symmetric matrix  $\mathbf{A}$  form an orthogonal basis such that any vector  $\mathbf{x}_0$  can be represented as a linear combination of the eigenvectors,

$$\mathbf{x}_0 = \sum_k c_k \mathbf{v}_k . \quad (10.3)$$

Acting on  $\mathbf{x}_0$  with the matrix  $\mathbf{A}$  gives

$$\mathbf{A}^i \mathbf{x}_0 = \sum_k \lambda_k^i c_k \mathbf{v}_k . \quad (10.4)$$

Thus the contribution from the eigenvector corresponding to the largest eigenvalue is amplified with the factor

$$\left( \frac{\lambda_{\text{largest}}}{\lambda_{\text{next largest}}} \right)^i . \quad (10.5)$$

The eigenvalue can be estimated using the *Rayleigh quotient*,

$$\lambda[\mathbf{x}_i] = \frac{\mathbf{x}_i^T \mathbf{A} \mathbf{x}_i}{\mathbf{x}_i^T \mathbf{x}_i} = \frac{\mathbf{x}_{i+1}^T \mathbf{x}_i}{\mathbf{x}_i^T \mathbf{x}_i}. \quad (10.6)$$

## 10.2 Inverse iteration

Alternatively, the *inverse power iteration* with the inverse matrix,

$$\mathbf{x}_{i+1} = \mathbf{A}^{-1} \mathbf{x}_i, \quad (10.7)$$

converges to the smallest (in the absolute value) eigenvalue of matrix  $A$ .

Finally, the *shifted inverse power iteration*,

$$\mathbf{x}_{i+1} = (\mathbf{A} - s\mathbf{1})^{-1} \mathbf{x}_i, \quad (10.8)$$

where  $\mathbf{1}$  signifies the identity matrix of the same size as  $\mathbf{A}$ , converges to the eigenvalue closest to the given number  $s$ .

The *inverse iteration method* is a refinement of the inverse power method where the trick is not to invert the matrix in (10.8) but rather solve the linear system

$$(\mathbf{A} - s\mathbf{1})\mathbf{x}_{i+1} = \mathbf{x}_i \quad (10.9)$$

using e.g. QR-decomposition.

The better approximation  $s$  to the sought eigenvalue is chosen, the faster convergence one gets. However, incorrect choice of  $s$  can lead to slow convergence or to the convergence to a different eigenvector. In practice the method is usually used when good approximation for the eigenvalue is known, and hence one needs only few (quite often just one) iteration.

One can update the estimate for the eigenvalue using the Rayleigh quotient  $\lambda[\mathbf{x}_i]$  after each iteration and get faster convergence for the price of  $O(n^3)$  operations per QR-decomposition; or one can instead make more iterations (with  $O(n^2)$  operations per iteration) using the same matrix  $(\mathbf{A} - s\mathbf{1})$ . The optimal strategy is probably an update after several iterations.

## 10.3 Krylov subspaces

When calculating an eigenvalue of a matrix  $\mathbf{A}$  using the power method, one starts with an initial random vector  $\mathbf{b}$  and then computes iteratively the sequence  $\mathbf{A}\mathbf{b}, \mathbf{A}^2\mathbf{b}, \dots, \mathbf{A}^{n-1}\mathbf{b}$  normalising and storing the result in  $\mathbf{b}$  on each iteration. The sequence converges to the eigenvector of the largest eigenvalue of  $\mathbf{A}$ .



The set of vectors

$$\mathcal{K}_n = \{\mathbf{b}, \mathbf{A}\mathbf{b}, \mathbf{A}^2\mathbf{b}, \dots, \mathbf{A}^{n-1}\mathbf{b}\}, \quad (10.10)$$

where  $n < \text{rank}(\mathbf{A})$ , is called the order- $n$  *Krylov matrix*, and the subspace spanned by these vectors is called the order- $n$  *Krylov subspace* [3]. The vectors are not orthogonal but can be made so e.g. by Gram-Schmidt orthogonalisation.

For the same reason that  $\mathbf{A}^{n-1}\mathbf{b}$  approximates the dominant eigenvector one can expect that the other orthogonalised vectors approximate the eigenvectors of the  $n$  largest eigenvalues.

Krylov subspaces are the basis of several successful iterative methods in numerical linear algebra, in particular: Arnoldi and Lanczos methods for finding one (or a few) eigenvalues of a matrix; and GMRES (Generalised Minimum RESidual) method for solving systems of linear equations.

These methods are particularly suitable for large sparse matrices as they avoid matrix-matrix operations but rather multiply vectors by matrices and work with the resulting vectors and matrices in Krylov subspaces of modest sizes.

## 10.4 Arnoldi iteration

Arnoldi iteration is an algorithm where the order- $n$  orthogonalised Krylov matrix  $\mathbf{Q}_n$  for a given matrix  $\mathbf{A}$  is built using stabilised Gram-Schmidt process:

```

start with a set  $\mathbf{Q} = \{\mathbf{q}_1\}$  where  $\mathbf{q}_1$  is a random normalised vector;
repeat for  $k = 2$  to  $n$  :
  make a new vector  $\mathbf{q}_k = \mathbf{A}\mathbf{q}_{k-1}$ 
  orthogonalise  $\mathbf{q}_k$  to all vectors  $\mathbf{q}_i \in \mathbf{Q}$  storing  $\mathbf{q}_i^\dagger \mathbf{q}_k \rightarrow h_{i,k-1}$ 
  normalise  $\mathbf{q}_k$  storing  $\|\mathbf{q}_k\| \rightarrow h_{k,k-1}$ 
  add  $\mathbf{q}_k$  to the set  $\mathbf{Q}$ 

```

By construction the matrix  $\mathbf{H}_n$  made of the elements  $h_{jk}$  is an upper Hessenberg matrix,

$$\mathbf{H}_n = \begin{bmatrix} h_{1,1} & h_{1,2} & h_{1,3} & \cdots & h_{1,n} \\ h_{2,1} & h_{2,2} & h_{2,3} & \cdots & h_{2,n} \\ 0 & h_{3,2} & h_{3,3} & \cdots & h_{3,n} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & h_{n,n-1} & h_{n,n} \end{bmatrix}, \quad (10.11)$$

which is a partial orthogonal reduction of  $A$  into Hessenberg form,

$$\mathbf{H}_n = \mathbf{Q}_n^\dagger \mathbf{A} \mathbf{Q}_n. \quad (10.12)$$

The matrix  $\mathbf{H}_n$  can be viewed as a representation of  $\mathbf{A}$  in the Krylov subspace  $\mathcal{K}_n$ . The eigenvalues and eigenvectors of the matrix  $\mathbf{H}_n$  approximate the largest eigenvalues of matrix  $\mathbf{A}$ .

Since  $\mathbf{H}_n$  is a Hessenberg matrix of modest size its eigenvalues can be relatively easily computed with standard algorithms.

In practice if the size  $n$  of the Krylov subspace becomes too large the method is restarted.

## 10.5 Lanczos iteration

Lanczos iteration is Arnoldi iteration for Hermitian matrices, in which case the Hessenberg matrix  $\mathbf{H}_n$  of Arnoldi method becomes a tridiagonal matrix  $\mathbf{T}_n$ .

The Lanczos algorithm thus reduces the original hermitian  $N \times N$  matrix  $\mathbf{A}$  into a smaller  $n \times n$  tridiagonal matrix  $\mathbf{T}_n$  by an orthogonal projection onto the order- $n$  Krylov subspace. The eigenvalues and eigenvectors of a tridiagonal matrix of a modest size can be easily found by e.g. the QR-diagonalisation method.

In practice the Lanczos method is not very stable due to round-off errors leading to quick loss of orthogonality. The eigenvalues of the resulting tridiagonal matrix may then not be a good approximation to the original matrix. Library implementations fight the stability issues by trying to prevent the loss of orthogonality and/or to recover the orthogonality after the basis is generated.

## 10.6 Generalised minimum residual (GMRES)

GMRES is an iterative method for the numerical solution of a system of linear equations,

$$\mathbf{Ax} = \mathbf{b}, \quad (10.13)$$

where the exact solution  $\mathbf{x}$  is approximated by the vector  $\mathbf{x}_n \in \mathcal{K}_n$  that minimises the residual  $\mathbf{Ax}_n - \mathbf{b}$  in the Krylov subspace  $\mathcal{K}_n$  of matrix  $\mathbf{A}$ ,

$$\mathbf{x} \approx \mathbf{x}_n \leftarrow \min_{\mathbf{x} \in \mathcal{K}_n} \|\mathbf{Ax} - \mathbf{b}\|. \quad (10.14)$$

The vector  $\mathbf{x}_n \in \mathcal{K}_n$  can be represented as  $\mathbf{x}_n = \mathbf{Q}_n \mathbf{y}_n$  where  $\mathbf{Q}_n$  is the projector on the space  $\mathcal{K}_n$  and  $\mathbf{y}_n$  is an  $n$ -dimensional vector. Substituting  $\mathbf{x}_n \in \mathcal{K}_n$  gives an overdetermined system

$$\mathbf{AQ}_n \mathbf{y}_n = \mathbf{b}, \quad (10.15)$$

which can be solved by the ordinary least-squares method.

One can also project equation (10.15) onto Krylov subspace  $\mathcal{K}_n$  which gives a square system

$$\mathbf{H}_n \mathbf{y}_n = \mathbf{Q}_n^\dagger \mathbf{b}, \quad (10.16)$$

where  $\mathbf{H}_n = \mathbf{Q}_n^\dagger \mathbf{A} \mathbf{Q}_n$ .



# Chapter 11

## Fast Fourier transform

Fast Fourier transform (FFT) is an efficient algorithm to compute the discrete Fourier transform (DFT).

Computing DFT of  $N$  points in the naïve way, using the definition, takes  $O(N^2)$  arithmetic operations, while an FFT can compute the same result in only  $O(N \log N)$  operations. The difference in speed can be substantial, especially for large data sets. This improvement made many DFT-based algorithms practical.

Since the inverse of a DFT is also a DFT, any FFT algorithm can be used in for the inverse DFT as well.

The most well known FFT algorithms, like the Cooley-Tukey algorithm [16], depend upon the factorization of  $N$ . However, there are FFTs with  $O(N \log N)$  complexity for all  $N$ , even for prime  $N$ .

### 11.1 Discrete Fourier Transform

For a set of complex numbers  $x_n$ ,  $n = 0, \dots, N - 1$ , the DFT is defined as a set of complex numbers  $c_k$ ,

$$c_k = \sum_{n=0}^{N-1} x_n e^{-2\pi i \frac{nk}{N}}, k = 0, \dots, N - 1. \quad (11.1)$$

The inverse DFT is given by

$$x_n = \frac{1}{N} \sum_{k=0}^{N-1} c_k e^{+2\pi i \frac{nk}{N}}. \quad (11.2)$$

These transformations can be viewed as expansion of the vector  $x_n$  in terms of the orthogonal basis of vectors  $e^{2\pi i \frac{kn}{N}}$ ,

$$\sum_{n=0}^{N-1} \left( e^{2\pi i \frac{kn}{N}} \right) \left( e^{-2\pi i \frac{k'n}{N}} \right) = N\delta_{kk'} . \quad (11.3)$$

The DFT represent the amplitude and phase of the different sinusoidal components in the input data  $x_n$ .

The DFT is widely used in different fields, like spectral analysis, data compression, solution of partial differential equations and others.

### 11.1.1 Applications

#### Data compression

Several lossy (that is, with certain loss of data) image and sound compression methods employ DFT as an approximation for the Fourier series. The signal is discretized and transformed, and then the Fourier coefficients of high/low frequencies, which are assumed to be unnoticeable, are discarded. The decompressor computes the inverse transform based on this reduced number of Fourier coefficients.

#### Partial differential equations

Discrete Fourier transforms are often used to solve partial differential equations, where the DFT is used as an approximation for the Fourier series (which is recovered in the limit of infinite  $N$ ). The advantage of this approach is that it expands the signal in complex exponentials  $e^{inx}$ , which are eigenfunctions of differentiation:  $\frac{d}{dx}e^{inx} = in e^{inx}$ . Thus, in the Fourier representation, differentiation is simply multiplication by  $in$ .

A linear differential equation with constant coefficients is transformed into an easily solvable algebraic equation. One then uses the inverse DFT to transform the result back into the ordinary spatial representation. This approach belongs to the group of *spectral methods*.

#### Convolution and Deconvolution

FFT can be used to efficiently compute convolutions of two sequences. A convolution is the pairwise product of elements from two different sequences, such as in multiplying two polynomials or multiplying two long integers.

Another example comes from data acquisition processes where the detector introduces certain (typically Gaussian) blurring to the sampled signal. A reconstruction of the original signal can be obtained by deconvoluting the acquired signal with the detector's blurring function.

## 11.2 Cooley-Tukey algorithm

In its simplest incarnation this algorithm re-expresses the DFT of size  $N = 2M$  in terms of two DFTs of size  $M$ ,

$$\begin{aligned}
 c_k &= \sum_{n=0}^{N-1} x_n e^{-2\pi i \frac{nk}{N}} \\
 &= \sum_{m=0}^{M-1} x_{2m} e^{-2\pi i \frac{mk}{M}} + e^{-2\pi i \frac{k}{N}} \sum_{m=0}^{M-1} x_{2m+1} e^{-2\pi i \frac{mk}{M}} \\
 &= \begin{cases} c_k^{(\text{even})} + e^{-2\pi i \frac{k}{N}} c_k^{(\text{odd})} & , k < M \\ c_{k-M}^{(\text{even})} - e^{-2\pi i \frac{k-M}{N}} c_{k-M}^{(\text{odd})} & , k \geq M \end{cases} \quad , \quad (11.4)
 \end{aligned}$$

where  $c^{(\text{even})}$  and  $c^{(\text{odd})}$  are the DFTs of the even- and odd-numbered sub-sets of  $x$ .

This re-expression of a size- $N$  DFT as two size- $\frac{N}{2}$  DFTs is sometimes called the Danielson-Lanczos lemma. The exponents  $e^{-2\pi i \frac{k}{N}}$  are called *twiddle factors*.

The operation count by application of the lemma is reduced from the original  $N^2$  down to  $2(N/2)^2 + N/2 = N^2/2 + N/2 < N^2$ .

For  $N = 2^p$  Danielson-Lanczos lemma can be applied recursively until the data sets are reduced to one datum each, see Table 11.1. The number of operations is then reduced to  $O(N \ln N)$  compared to the original  $O(N^2)$ .

The established library FFT routines, like FFTW and GSL, further reduce the operation count (by a constant factor) using advanced programming techniques like precomputing the twiddle factors, effective memory management and others.

## 11.3 Multidimensional DFT

For example, a two-dimensional set of data  $x_{n_1 n_2}$ ,  $n_1 = 1 \dots N_1$ ,  $n_2 = 1 \dots N_2$  has the discrete Fourier transform

$$c_{k_1 k_2} = \sum_{n_1=0}^{N_1-1} \sum_{n_2=0}^{N_2-1} x_{n_1 n_2} e^{-2\pi i \frac{n_1 k_1}{N_1}} e^{-2\pi i \frac{n_2 k_2}{N_2}} \quad . \quad (11.5)$$

Table 11.1: C-implementation of the Cooley-Tukey algorithm

```

#include<complex.h>
#include<tgmath.h>

void dfts(int sign, int N, complex* x, int stride, complex* c){
    for(int k=0;k<N;k++){
        c[k]=0;
        for(int n=0;n<N;n++) c[k]+=x[n*stride]*exp(sign*2*M_PI*I*n*k/N);
    }
}

void fft2s(int sign, int N, complex*x, int stride, complex*c){
    if(N==1) c[0]=x[0];
    else if(N%2==0){
        fft2s(sign,N/2,x,2*stride,c);
        fft2s(sign,N/2,x+stride,2*stride,c+N/2);
        for(int k=0;k<N/2;k++){
            complex p=c[k], q=exp(sign*2*M_PI*I*k/N)*c[k+N/2];
            c[k]=p+q;
            c[k+N/2]=p-q;
        }
    }
    else dfts(sign,N,x, stride,c);
}

void fft(int N,complex*x,complex*c){ fft2s(-1,N,x,1,c); }
void ift(int N,complex*c,complex*x){
    fft2s(+1,N,c,1,x);
    for(int i=0;i<N;i++)x[i]/=N;
}

```



# Bibliography

- [1]
- [2] Hiroshi Akima. A new method of interpolation and smooth curve fitting based on local procedures. *Journal of Assoc. for Comp. Mach.*, 17(4):589–602, 1970.
- [3] A.N.Krylov. On the numerical solution of equation by which are determined in technical problems the frequencies of small vibrations of material systems. *Izvestiia Akademii nauk SSSR (in Russian)*, 7(4):491–539, 1931.
- [4] R.M. Barnett. Review of particle properties. *Physical Review D*, 54:1, 1996.
- [5] J.P. Berrut. Rational functions for guaranteed and experimentally well-conditioned global interpolation. *Comput. Math. Appl.*, 15(1):1–16, 1988.
- [6] Przemyslaw Bogacki and Lawrence F. Shampine. A 3(2) pair of Runge–Kutta formulas. *Applied Mathematics Letters*, 2(4):321–325, 1989.
- [7] C.G. Broyden. A class of methods for solving nonlinear simultaneous equations. *JSTOR*, 19(92):577–593, 1965.
- [8] C.W. Clenshaw and A.R. Curtis. A method for numerical integration on an automatic computer. *Numerische Mathematik*, 2:197–205, 1960.
- [9] M. Galassi et al. *GNU Scientific Library Reference Manual*. Network Theory Ltd, 3rd edition, 2009.
- [10] Erwin Fehlberg. *Low-order classical Runge-Kutta formulas with step size control and their application to some heat transfer problems*. NASA Technical Report, 1969.
- [11] Wallace Givens. Computation of plane unitary rotations transforming a general matrix to triangular form. *J. SIAM*, 6(1):26–50, 1958.

- [12] Gene H. Golub. Some modified matrix eigenvalue problems. *SIAM Rev.*, 15(2):318–334, 1973.
- [13] Gene H. Golub and John H. Welsch. Calculation of Gauss quadrature rules. *Mathematics of Computation*, 23(106):221–230, 1969.
- [14] Pedro Gonnet. Increasing the reliability of adaptive quadrature using explicit interpolants. *ACM Trans. Math. Soft.*, 37(3):26:2–26:32, 2010.
- [15] A.S. Householder. Unitary triangularization of a nonsymmetric matrix. *Journal of the ACM*, 5(4):339–342, 1958.
- [16] John.W.Tukey James.W.Cooley. An algorithm for the machine calculation of complex fourier series. *Math. Comput.*, 19(90):297–301, 1965.
- [17] J.A.Nelder and R.Mead. A simplex method for function minimization. *Computer Journal*, 7:308–313, 1965.
- [18] Aleksandr Semenovich Kronrod. *Nodes and weights of quadrature formulas. Sixteen-place tables.* Consultants Bureau, 1965.
- [19] Masatake Mori. Quadrature formulas obtained by variable transformation and the DE–rule. *Journal of Computational and Applied Mathematics*, 12&13:119–130, 1985.
- [20] T. N. L. Patterson. The optimum addition of points to quadrature formulae. *Mathematics of Computation*, 22(104):847–856, 1968.
- [21] Robert Piessens, Elise de Doncker-Kapenga, Christoph W. Überhuber, and David Kahaner. *QUADPACK: A subroutine package for automatic integration.* Springer-Verlag, 1983.
- [22] B. Riemann. Über die darstellbarkeit einer function durch eine trigonometrische reihe. *Abhandlungen der Königlichen Gesellschaft der Wissenschaften zu Göttingen*, 13:87–132, 1868.