Excercises in Numerical Methods



Fredrik Berntsson (fredrik.berntsson@liu.se)

Linköpings universitet

Contents

1	Basic Concepts and Floating Point systems	3
2	Error Analysis	5
3	Non-linear equations	7
4	Polynomial and Spline Interpolation	10
5	Integration, Differentiation and Extrapolation	15
6	Differential Equations	18
7	Basic matrix operations and linear systems of equations	19
8	Least Squares Problems and Orthogonal Decompositions	23
A	nswers	26
1	Basic Concepts and Floating Point systems	26
2	Error Analysis	28
3	Non-linear equations	31
4	Polynomial and Spline Interpolation	35
5	Integration, Differentiation and Extrapolation	39
6	Differential Equations	42
7	Basic matrix operations and linear systems of equations	44
8	Least Squares Problems and Orthogonal Decompositions	49

1 Basic Concepts and Floating Point systems

Exercise 1.1 Let $a = 0.0987 \pm 0.5 \cdot 10^{-4}$ and $b = 20.104 \pm 4 \cdot 10^{-3}$. Determine the number of correct decimals and the number of significant digits for both a and b.

Exercise 1.2 Let $c_0 \approx \bar{c}_0 = 2.99792458 \cdot 10^6$ be correctly rounded. How many correct decimals and significant digits does the approximate value \bar{c}_0 have?

Exercise 1.3 Let a = 22.73531443. Round the value *a* correctly to 5 significant digits to obtain the approximation \bar{a} . Give both the approximate value \bar{a} and an upper bound for the absolute error $|\Delta a|$ in the approximation.

Exercise 1.4 We approximate π by $\bar{\pi} = 3.1415$. How many correct decimals and significant digits do we have?

Hint $\pi = 3.1415926535...$

Exercise 1.5 Round the mathematical constant of the golden ratio $\phi = 1.61803398875...$ to five correct decimals and also to five significant digits.

Exercise 1.6 Let $\bar{a} = 22.73531$ be an approximate value of a. The error in \bar{a} is $|\Delta a| \le 4.7 \cdot 10^{-3}$. How many *correct decimlals* and *significant digits* does the approximation \bar{a} have?

Exercise 1.7 Let $y = e^{\pi}$. Classify the following error sources in the computation of y: The rounding of π to 3.142 and the approximation $e^x \approx 1 + x + \frac{1}{2}x^2$.

Exercise 1.8 Let x = -102.232. Give a bound for the *absolute error* when x is stored on a computer using the floating point system (10, 3, -10, 10).

Exercise 1.9 Let x = 168.3556541. What is the closest number to x that exists in the floating point system (10, 5, -100, 100).

Exercise 1.10 A single precision floating point number $x = (-1)^s (1.f)_2 2^{e-127}$ is stored using 32 bits assigned as follows

Clearly show how the numbers x = 9.25 and x = -2.65625 are stored.

Exercise 1.11 Consider the floating point number system (10, 5, -9, 9), where $\beta = 10$ is the base and t = 5 is the number of digits in the fractional part. Let x = 117.5614 and y = 0.01678214. Find the closest numbers x_r and y_r in the floating point system.

Exercise 1.12 Suppose we represent a number x in the form $x = m \times \beta^e$. What do we call m, β and e? Also, What else must be specified in order to properly define a floating point number system?

Exercise 1.13 Consider the floating point system (10, 2, -9, 9). Let $a = 8.50 \cdot 10^5$ and $b = 5.25 \cdot 10^2$. Compute the floating point results fl[$a \cdot b$] and fl[a/b]. In both cases also give a bound for the relative error in the result.

Exercise 1.14 Rewrite the expressions $\sqrt{1+x} - 1$, $(1-x)^{-1} - (1+x)^{-1}$, and $1 - \cos^2(x)$ in such a way that the cancellation is avoided.

Exercise 1.15 Let $\sqrt{101} = 10.04988$ be correctly rounded so that the error is at most $0.5 \cdot 10^{-5}$. What is the resulting absolute error in $x = \sqrt{101} - 10$? Also propose an alternative formula for computing x that leads to a smaller error.

2 Error Analysis

Exercise 2.1 Let f = (x - y)/z, where x = 8.25, y = 1.05 and z = 4.00 are correctly rounded. Compute an approximate value for f together with a bound for the absolute error.

Exercise 2.2 Let $f(x) = e^{bx}$, where $b = 1.70 \pm 0.01$. Compute f(2) and also use the error propagation formula to give an error bound.

Exercise 2.3 The area of a circle is $A = \pi r^2$. Compute the area, and an error bound, for the case when $r = 23.76 \pm 0.02$ and we approximate π by 3.142.

Exercise 2.4 Let $y = ae^b$, where $a = 1.54 \pm 0.03$ and $b = 3.17 \pm 0.05$. Compute the approximate value \bar{y} and an error bound.

Exercise 2.5 The focal point of a lens can be determined by the formula

$$\frac{1}{f} = \frac{1}{a} + \frac{1}{b}$$
, where $a = 32 \pm 1$ and $b = 46 \pm 1.3$

Determine f(a, b) and an error bound.

Exercise 2.6 We compute the function

$$f(x) = \sqrt{1+x} - \sqrt{1-x}$$

for small x values on a computer with unit round off $\mu = 1.11 \cdot 10^{-16}$. We find that the results are quite poor and that the *relative error* in the result tends to grow as $x \to 0$. Explain the poor accuracy by performing an analysis of the computational errors and give a bound for the relative error in the computed result f(x). For the analysis you may assume that all computations are performed with a relative error at most μ .

Exercise 2.7 We want to evaluate a function f(x) on a computer, for small values of x, and have two alternate expressions:

$$f_1(x) = \frac{1 - \cos(x)}{\sin(x)}$$
 or $f_2(x) = \frac{\sin(x)}{1 + \cos(x)}$

For the case $x = 1.111 \cdot 10^{-8}$ we evaluate both expressions in Matlab and obtain $f_1 = 9.9930 \dots 10^{-9}$ and $f_2 = 5.5550 \dots 10^{-9}$.

Assume that all numerical computations are performed with a relative error of at most the unit round off μ and perform an analysis of the computational errors. Derive a bound for both the absolute error in the results.

Hint The unit round off for Matlab is $\mu = 1.11 \cdot 10^{-16}$.

Exercise 2.8 We compute the function

$$f(x) = e^x - 3x$$

for small x values on a computer with unit round off $\mu = 1.11 \cdot 10^{-16}$. Preform an analysis of the computational errors to obtain a bound for the relative error in the computed results f(x). For the analysis you may assume that all computations are performed with a relative error at most μ . Also, use the obtained bound to argue if *cancellation* occurs during the computations. In case of cancellation also suggest an alternative formula that can be expected to give better accuracy.

Exercise 2.9 We want to evaluate the function

$$f(x) = \frac{x - \sin(x)}{x^3}$$

for small values of x on a computer with the unit round off $\mu = 1.11 \cdot 10^{-16}$. A Taylor series expansion shows that

$$\lim_{x \to 0} f(x) = \frac{1}{6},$$

but when we compute the expression for $x = 10^{-7}$ we get the difference $|\bar{f}(10^{-7}) - \frac{1}{6}| \approx 5.4 \cdot 10^{-3}$. Explain the above result by performing an error analysis that clearly shows how large the error is when f(x) is evaluated on the computer. For the analysis you should assume that all computations are carried out with a relative error of at most μ .

Hits The Taylor series expansion of sin(x) is $sin(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots$

Exercise 2.10 We compute the function

$$f(x) = 1 - 2x\cos(x)$$

for small x values on a computer with unit round off $\mu = 1.11 \cdot 10^{-16}$. Preform an analysis of the computational errors to obtain a bound for the relative error in the computed results f(x). For the analysis you may assume that all computations are performed with a relative error at most μ . Also, use the obtained bound to argue if *cancellation* occurs during the computations. In case of cancellation also suggest an alternative formula that can be expected to give better accuracy.

Exercise 2.11 Assume that we have an approximate value \bar{x} and want to find the resulting error Δf . The general error propagation formula states that

$$|\Delta f| \lesssim |\frac{\partial f}{\partial x}| |\Delta x|.$$

If $|f'(\bar{x})| \approx |f'(x)| = 0$ the formula fails. Show that for this case it is more resonable to use

$$|\Delta f| \lesssim |f''(\bar{x})| \frac{|\Delta x|^2}{2}.$$

3 Non-linear equations

Exercise 3.1 Show that the function $f(x) = e^x - \frac{4}{2+x}$ has a root in the interval [0, 1].

Exercise 3.2 The equation $e^x = 2(1 - x)$ has a root $x^* \approx 0.3$. Does the fixed point iteration $x_{k+1} = 1 - \frac{1}{2}e^{x_k}$ converge to the root?

Exercise 3.3 We have solved the equation $f(x) = x^3 + x - 7 = 0$ and obtained an approximate root $\bar{x} = 1.7$. Estimate the error in the approximation \bar{x} .

Exercise 3.4 We have solved f(x) = 0 using Newtons method and obtained an approximate root \bar{x} such that $f(\bar{x}) = 0$ when we evaluate the function on the computer. Suppose we can compute f(x) with an absolute error at most 10^{-8} and that $1.7 \leq f'(x) \leq 2.2$ near the root x^* . Estimate the error in \bar{x} for this case.

Exercise 3.5 Determine the positive root of the equation $x = 5(1 - e^{-x})$ with five correct decimals.

Exercise 3.6 We are intressted in solving the equation $x = 2\sin(x)$. Two possible fixed point iterations are obtained by using the iteration functions $\phi_1(x) = 2\sin(x)$ or $\phi_2(x) = \frac{x}{2} + \sin(x)$. Which of these fixed point iterations would have the fastest convergence if $x^* \approx 1.9$?

Exercise 3.7 Consider the equation $f(x) = 2\cos(x) - 4x$ which has a positive root $x^* \approx 0.45$. Show theoretically that the fixed point iteration $x_{k+1} = \phi(x_n) = \frac{1}{2}\cos(x_n)$ is convergent for any starting value x_0 . Also, if we use $x_0 = 0.45$ then we get

$$\begin{array}{c|ccc} k & x_k \\ \hline 0 & 0.4502236 \\ 1 & 0.4501749 \\ 2 & 0.4501855 \\ 3 & 0.4501832 \end{array}$$

Estimate the error in the approximation $\bar{x} = 0.4501832$.

Exercise 3.8 The equation $f(x) = x - 3e^{-x} = 0$ has a solution $x^* \approx 1.05$.

- a) Estimate the error in the approximation $\bar{x} = 1.05$.
- **b**) Investigate the following fixed point iterations theoretically

(i)
$$x_{n+1} = 3e^{-x_n}$$
, (ii) $x_{n+1} = (2x_n + 3e^{-x_n})/3$,
(iii) $x_{n+1} = 1.05x_n + 3e^{-x_n}$, (iv) $x_{n+1} = (x_n + 3e^{-x_n})/2$

Determine if they converge towards x^* and also find out which of the methods that have the fastest convergence.

c) Estimate the number of iterations that would be needed for the fastest of the methods from b) if $x_0 = 1.05$ and we want to find the root with an absolute error of at most 10^{-10} .

Exercise 3.9 Show that Newtons method is convergent when applied to the equation $f(x) = x^2 = 0$ and if the starting guess is $x_0 = 1$. Also show that the order of convergence is p = 1.

Exercise 3.10 Show that Newton-Raphsons has a quadratic rate of convergence if x^* is a single root.

Exercise 3.11 The equation $f(x) = x^3 - 7.5x^2 + 18x - 14 = 0$ has a double root x^* . Using the secant method we have obtained $\bar{x} = 1.99789$. Derive the error estimate

$$|\bar{x} - x^*|^2 = 2 \left| \frac{f(\bar{x})}{f''(\xi)} \right|, \quad \xi \in (\bar{x}, x^*),$$

which is valid for double roots. Also compute the error bound for the approximate root $\bar{x} = 1.99789$.

Exercise 3.12 Consider the equation $f(x) = \cos(x) - xe^x = 0$. We use the Newton-Raphson method for finding the root, with the starting guess $x_0 = 1$, and obtain the following table

k	x_k	$f(x_k)$
0	1.0000000	$-2.2 \cdot 10^{0}$
1	0.6530794	$-4.6 \cdot 10^{-1}$
2	0.5313434	$-4.2\cdot10^{-2}$
3	0.5179099	$-4.6\cdot10^{-4}$
4	0.5177574	$-5.9\cdot10^{-8}$

We decide to use $\bar{x} = 0.5178$ as an approximation of x^* . Estimate the error in the approximation \bar{x} . Also state the definition of the *order of convergence* for an iterative method and use the table above to estimate the order of convergence for the Newton-Raphson method.

Exercise 3.13 The Newton-Raphsson method is used to find a root of the cubic polynomial $f(x) = x^3 - 9x^2 + 24x - 20$. We obtain the following iteration sequence

k	x_k	$ x_k - x^* $
0	1.8000	0.2000
1	1.8970	0.1030
2	1.9476	0.0524
3	1.9736	0.0264
4	1.9867	0.0133

State the definition of order of convergence p for an iterative method. Also use the table to determine the order of convergence when Newton-Raphson's method is applied to this specific function f(x). Also, use the results and known properties of Newton-Raphson's method to determine if $x^* = 2$ is a double or single root.

Exercise 3.14 We wish to implement the standard function $\sqrt{\cdot}$ on a computer using double precision arithmetic. We compute $x = \sqrt{a}$ by solving the equation $f(x) = x^2 - a$.

- a) Derive the iteration formula $x_k = \varphi(x_{k-1})$ obtained by applying Newton-Raphsons method to the above equation.
- **b)** Explain clearly why it is sufficient to consider the case $1 \le a < 4$, and thus $1 \le x < 2$.
- c) A convergence analysis for the Newton-Raphson method leads to the estimate,

$$|x_k - \sqrt{a}| \lesssim \frac{1}{2} |\varphi''(\sqrt{a})| |x_{k-1} - \sqrt{a}|^2.$$

Take advantage of this and determine the number of iterations needed to achive an error bound $|x_k - \sqrt{a}| \le \mu = 1.1 \cdot 10^{-16}$ if $x_0 = 1.5$ is used.

d) We wish to decrease the number of iterations, by picking a better starting guess, and select evenly spaced numbers $a_j = 1 + 3j/n$, j = 0, ..., n - 1, and compute $\sqrt{a_j}$ exactly. Clearly demonstrate how the values $(\sqrt{a_j}, a_j)$ can be used to obtain a better starting guess. What table size n do we need to reduce the number of iterations by one compared to the result in **c**)?

4 Polynomial and Spline Interpolation

Exercise 4.1 Let $p_n(x)$ be a polynomial of degree n. How many interpolation conditions of the type $p_n(x_i) = f_i$ are needed for $p_n(x)$ to be uniquely determined?

Exercise 4.2 Let the following table with correctly rounded function values be given

x	0.9	1.1	1.2
f(x)	0.4710	0.2452	0.2385

Find an approximate value for f(1.03). Also provide a complete error estimate.

Exercise 4.3 In an application we need to implement the function $y = \log(x)$, for $1 \le x \le 4$, with a maximum error $\varepsilon \le 10^{-5}$. We decide to use linear interpolation and create a table $\{x_i, y_i\}_{i=1}^n$, where $x_1 = 1$, $x_n = 4$ and $h = x_{i+1} - x_i$ is the stepsize. In the table we store approximate values $y_i \approx \log(x_i)$, rounded to 6 correct digits. Determine the smallest size n for the table so that the maximum error in the interpolated values is less than 10^{-5} .

Exercise 4.4 Let p(x) be the linear polynomial that interpolates the function $f(x) = \sin(x)$ at the points x = 0 and x = 1. Show that the truncation error is bounded by $|p(x) - \sin(x)| \le \frac{1}{8}$, for 0 < x < 1.

Exercise 4.5 A table with correctly rounded function values is given.

x	0.0	0.5	1.0	1.5
f(x)	1.80	2.80	4.10	5.90

Use linear interpolation to find an approximation of f(0.4) and also give an error estimate.

Exercise 4.6 The following table is given

x	0.6	0.7	0.8	0.9
f(x)	1.23	1.29	1.32	1.07

Use quadratic interpolation and compute an approximate value for f(0.74). Also estimate the truncation error in the result.

Exercise 4.7 Let x_1 , x_2 , x_3 and x_4 be given interpolation points. In the Lagrange interpolation formula we use basis functions $\ell_i(x)$ such that $\ell_i(x_j) = 1$ if i = j and zero otherwise. Give an explicit expression for the basis function $\ell_2(x)$ for the case with n = 4 interpolation points. What is the degree of the basis polynomial?

Exercise 4.8 Use Lagrange interpolation to find the polynomial of degree 2 that interpolates the table

x	1	2	3
f(x)	1.3	0.6	1.9

Exercise 4.9 Let $p(x) = c_0 + c_1 x + c_2 x^2 + c_3 x^3$ be a cubic polynomial. We want to find values for the coefficients so that p(0) = p(1) = 0 and p'(0) = p'(1) = 1. Show how to derive a linear system of equations such that the solution $c = (c_0, c_1, c_2, c_3)^T$ are the coefficients of a cubic polynomial satisfying these conditions. Also find the specific polynomial satisfying all the above conditions.

Exercise 4.10 Spline interpolation can be used to approximate a function y = f(x). We have a table

x	-2	-1	0	1	2
f(x)	0	1	3	1	0

We attempt to approximate f(x) by a cubic spline s(x). Clearly state the conditions that have to be satisfied for s(x) to be a cubic spline that interpolates the above table. Is the given information sufficient for the spline s(x) to be uniquely determined?

Exercise 4.11 Let

$$s(x) = \begin{cases} x+1 & 0 \le x < 1, \\ x^3 - 3x^2 + 4x & 1 \le x < 2. \end{cases}$$

Is s(x) a cubic spline?

Exercise 4.12 Let

$$s(x) = \begin{cases} ax + 1 & 0 \le x < 1, \\ bx^3 + cx^2 & 1 \le x < 2. \end{cases}$$

Determine the constants a, b and c so that s(x) is a cubic spline.

Exercise 4.13 Approximate the function $f(x) = x^3 + x^2 + 1$ by a cubic spline s(x) that interpolates f(x) at the nodes x = 0, 0.3, 0.6, 0.7 and 1.0. Use correct end point conditions, i.e. s'(0) = f'(0) and s'(1) = f'(1). Give the expression for s(x).

Exercise 4.14 Consider a case where s(x) is defined by two cubic polynomials,

$$s(x) = \begin{cases} s_1(x) = 0.9 + 0.1x + 0.6x^2 + ax^3, & 0 \le x < 1, \\ s_2(x) = 2.0 + b(x-1) + c(x-1)^2 + 0.4(x-1)^3, & 1 \le x \le 2. \end{cases}$$

Find the appropriate values for the constants a, b and c so that s(x) is a cubic spline that interpolates the table tabellen:

х	0	1.0	2
s(x)	0.9	2.0	6.7

with the end point conditions s'(0) = 0.1 and s'(2) = 7.3.

Exercise 4.15 A function s(x) is given by two cubic polynomials

$$s(x) = \begin{cases} s_1(x) = 0.8 + 0.2x - 0.4x^2, & 0 \le x < 1, \\ s_2(x) = 0.6 - 0.6(x - 1) - 0.4(x - 1)^2 - 0.4(x - 1)^3, & 1 \le x \le 2. \end{cases}$$

Is s(x) a cubic spline? Present the calculations. Also determine if s(x) is a natural cubic spline?

Exercise 4.16 A function f(x) can be approximated by a piecewise polynomial s(x) on the interval [a, b] by introducing evenly spaced nodes

$$a = x_0 < x_1 < x_2 < \ldots < x_N = b$$

On each subinterval $[x_k, x_{k+1}]$ we let s(x) be given by a cubic polynomial

$$s_k(x) = a_k + b_k(x - x_k) + c_k(x - x_k)^2 + d_k(x - x_k)^3, \qquad x_k \le x < x_{k+1}.$$

Clearly formulate the conditions that needs to be satisfied for s(x) to be a cubic spline, defined on [a, b], and that interpolates f(x) in the nodes $\{x_k\}_{k=0}^N$.

Also illustrate the case N = 3 and draw a sketch that clearly illustrates the nodes, interpolation points and polynomials.

Exercise 4.17 In order to obtain a unique interpolating spline s(x) we use correct end point conditions, i.e. s'(a) = f'(a) and s'(b) = f'(b). We experiment with different number of nodes N and measure the maximum error max |s(x) - f(x)| on the interval. This gives us the table

N	5	10	20
$\max s(x) - f(x) $	0.0435	0.00269	0.000172

We know that the maximum error should depend on the step size $h = \max |x_{k+1} - x_k|$ as Ch^p , where p is an integer and C is a constant. Use the numbers in the table to determine p.

Exercise 4.18 Construct a linear spline that interpolates the table

x	1	3	4
f(x)	1.56	2.31	1.97

Exercise 4.19 A quadratic Beziér curve is given by the expression

$$p(t) = (1-t)^2 P_1 + 2(1-t)t P_2 + t^2 P_3, \quad 0 < t < 1,$$

where P_1 , P_2 and P_3 are control points.

- a) Show that the tangent, for t = 0 is parallell to the vector $P_2 P_1$.
- **b**) Suppose we want to put together two quadratic Beziér curves. We select five control points according to the sketch:



The point P_3 is common for both curve segments. We have $P_2 = (2, 6)^T$, $P_3 = (3, 5)^T$ and $P_5 = (6, 1)$. Clearly show how to pick the point P_4 so that the tangent direction of the curve is continuous in at P_3 and so that the tangent direction at P_5 is vertical.

Exercise 4.20 $P_1 = (1, 0)^T$, $P_2 = (1, 3)^T$, $P_3 = (4, 3)^T$ and $P_4 = (4, 2)^T$. Draw a sketch that clearly shows the convex hull formed by these points. Also use the available information to draw the cubic Beziér curve formed by the four points P_1, \ldots, P_4 as accurately as possible.

Exercise 4.21 A cubic Beziér curve is given by

$$p(t) = (1-t)^3 P_1 + 3(1-t)^2 t P_2 + 3(1-t)t^2 P_3 + t^3 P_4, \quad 0 < t < 1,$$

where P_1 , P_2 , P_3 and P_4 are control points.

- a) Show that the tangent of the curve in the starting points t = 0 is parallel to the vector $P_2 P_1$.
- b) Give the definition of the *convex hull* formed by the points P_1 , P_2 , P_3 and P_4 . Also show that the cubic Beziér curve is located within the convex hull formed by its control points P_1 , P_2 , P_3 and P_4 .
- c) Let $P_1 = (0, 0)^T$, $P_2 = (1, 3)^T$, $P_3 = (4, 2)^T$ and $P_4 = (5, 1)^T$. Also let s(t) be the cubic Beziér curve given by these control points. Compute s(1/2) and use the available information to draw a sketch that, as accurately as possible, shows the shape of the curve s(t).

Exercise 4.22 Create a parametric curve composed of two cubic Beziér curves as shown in the figure



Chose the points P_3 and P_5 so that the curve has a continuous tanget direction at P_4 . Also make sure the slope is exactly -2, e.g. the line from P_3 to P_5 can be expressed as y = -2x + b for some constant b.

5 Integration, Differentiation and Extrapolation

Exercise 5.1 Suppose a function f(x) have the required number of continuous derivatives. We want to approximate the derivative f'(x) by the difference formula

$$Df(x) = \frac{1}{h}(f(x+h) - f(x)).$$

Show that the truncation error can be written as $f'(x) - Df(x) \approx Ch$. Give an expression for C.

Exercise 5.2 The the function f(x) be known at the points $\{x_k\}$ by the table

x	0.0	0.25	0.5	0.75	1.0
f(x)	1.000	1.015	1.006	0.882	0.670

with correctly rounded function values. Use the Trapezoidal method to compute an approximation of the integral

$$I = \int_0^1 f(x) dx,$$

using the stepsize h = 0.25. Estimate the error in the result.

Exercise 5.3 Suppose a function f(x) can be computed with a *relative error* at most ε . Derive a bound for the resulting error when we want to compute an integral

$$I = \int_{a}^{b} f(x) dx,$$

using Simpsons rule. Also discuss if there are any cases where the error bound indicates that the the computation may be problematic.

Exercise 5.4 To compute the derivative f'(2) we can use the formula

$$Df(2) = \frac{1}{2h}(-f(x+2h) + 4f(x+h) - 3f(x)).$$

When the formula is applied for a few different h values we obtain the results

Assume that the error is proportional to h^p and use the table to determine p.

Exercise 5.5 Suppose $F_1(h) = a + bh^{p_1} + ch^{p_2}$, for some constants a, b and c and positive integers p_1 and p_2 such that $p_1 < p_2$. Show that

$$F_1(h) + \frac{F_1(h) - F_1(qh)}{q^{p_1} - 1} = a + \mathcal{O}(h^{p_2}).$$

Exercise 5.6 We compute an approximation of an integral using a numerical method T(h) with a truncation error $R_T \approx Ch^4$, where C is a constant. We use two different values for h and obtain

Use the table to estimate the truncation error in the computed value T(0.05).

Exercise 5.7 A numerical method has a truncation error that can be written as $R_T = Ch^p$, where C is a constant, p > 0 is an integer, and h is a discretization parameter. The method computes a value T(h) which approximates the exact value T_0 . We compute T(h) for a few different h-values to obtain the table

Use the table to determine p.

Exercise 5.8 Compute an approximation of $\int_0^1 f(x) dx$ using the following table with correctly rounded function values

by the Trapezoidal method. Use the step size h = 1/4. Also estimate the error in the approximate value.

Exercise 5.9 We use a numerical method to compute the derivative of a function f(x) and obtain

h	0.1	0.05	0.025
derivative	0.69280	0.71195	0.72171
error	-0.0388	-0.0196	-0.0099

The dominating source of error is the truncation error R_T which can be assumed to depend on h as $R_T \approx Ch^p$. Use the table to compute both C and p. Also determine the largest possible step size that can be used if we require that $|R_T| \leq 10^{-8}$.

Exercise 5.10 We want to compute an approximation of

$$\int_{1}^{4} f(x)dx, \quad \text{where } f(x) = \cos(x^{2})\sqrt{4-x}.$$

Do you expect the Trapezoidal method to work well for this case? Explain your conclusion clearly.

Exercise 5.11 We want to compute the value of the integral

$$I = \int_0^1 \mathrm{e}^{1-x^2} dx,$$

using Simpsons rule. We use a few different step sizes h and obtain

- (a) Suppose $S(h) = I + R_T$, where $R_T \approx c \cdot h^p$ for constants c and p. Clearly demonstrate how the table can be used to determine values for c and p. Answer with both a formula and the resulting values for c and p.
- (b) Suppose we want to find an approximate value for the integral with a total error $|R_T| < 10^{-6}$. What step size would be required? Motivate your answer.

6 Differential Equations

Exercise 6.1 Consider the problem

$$\begin{cases} y' = y^2 - t. \\ y(0) = 1, \end{cases}$$

Use the Euler method, and the stepsize h = 0.5, to compute an approximation of y(0.5). With the stepsize h = 0.1 the Euler method gives an approximation $y(0.5; h = 0.1) \approx 1.668$. Use this information to estimate the truncation error in the approximation y(0.5; h = 0.5).

Exercise 6.2 Consider the ordinary differential equation $y' = t^2 - y$, with y(0) = 2. Compute approximations of the solution y(0.2) using the Euler method and stepsizes h = 0.2 and h = 0.1. Also estimate the error in the approximation y(0.2; h = 0.1).

Exercise 6.3 Consider the problem $y' = -100y + \sin^2(t)$, y(0) = 0. Determine the maximum stepsize h such that the Euler method is stable when applied to the problem.

Exercise 6.4 We want to solve the problem y' = -25y + t + 0.04, y(0) = 1, using the Euler method. We require that thre solution is stable. Find the largest possible time step h that can be used.

Exercise 6.5 The Trapezoidal method can be written

$$y_{k+1} = y_k + \frac{h}{2}(f(t_k, y_k) + f(t_{k+1}, y_{k+1})).$$

Determine if the method is explicit or implicit. Motivate your answer. Also show that the method is stable for $h\lambda < 0$. This means that the trapezoidal method is always stable.

Exercise 6.6 Heun's method computes $y_k \approx y(t_k)$ by the steps

$$k_1 = hf(t_k, y_k), k_2 = hf(t_k + h, y_k + k_1), y_{k+1} = y_k + \frac{1}{2}(k_1 + k_2).$$

Derive a condition that guarantees that Heun's method is stable when applied to the test problem $y' = \lambda y, y(0) = 1.$

Exercise 6.7 Rewrite the initial value problem y'' + 3y' - ty + 1 = 0, y(0) = 1, y'(0) = 0, as a system of first order.

Exercise 6.8 We have the second order equation

$$y'' = 2y(1+y^2), \quad y(0) = 1, \ y'(0) = 2.$$

Rewrite the problem as a system of first order and compute an approximation of y(0.4) using the Euler method and h = 0.2.

Exercise 6.9 Van der Pools equation is

$$y'' + e(y^2 - a)y' + y = 0, \quad y(0) = a, \ y'(0) = b.$$

Rewrite the equation as a system of first order equations.

7 Basic matrix operations and linear systems of equations

Exercise 7.1 Suppose $A \in \mathbb{R}^{m \times m}$ and $B \in \mathbb{R}^{m \times n}$, m > n. How many operations are required to evaluate the formula z = (A + I)Bx + y, where x and y are vectors.

Exercise 7.2 Suppose A, B, and C are matrices and b is a vector. How would you implement the formula

$$x = B^{-1}(2A + I)(C^{-1} + A)b.$$

without computing any matrix inverse? Aim for as few arithmetic operations as possible.

Exercise 7.3 Suppose we want to solve the upper triangular system Rx = y by backwards substitution. Clearly show how many floating point operations are needed.

Exercise 7.4 The following table shows the time t, in seconds, needed to solve linear systems of equations with n unknowns on a computer.

The *LU*-decomposition followed by two triangular systems is used. Is the algorithms computational complexity as expected?

Exercise 7.5 The unit circle can be defined as all points $x = (x_1, x_2)$ such that ||x|| = 1. Draw the unit circle when $|| \cdot ||_1$, $|| \cdot ||_2$ and $|| \cdot ||_{\infty}$ is used.

Exercise 7.6 Explain what is ment by a matrix norm beeing *induced* from a vector norm. Also show that if A and B are matrices then for an induced norm $||AB|| \leq ||A|| ||B||$.

Exercise 7.7 Prove the inequality $||x||_{\infty} \leq ||x||_2 \leq \sqrt{n} ||x||_{\infty}$.

Exercise 7.8 Let $u \in \mathbb{R}^m$ and $v \in \mathbb{R}^n$. Show that

$$||uv^T||_2 = ||u||_2 ||v||_2.$$

Exercise 7.9 Prove that ||I|| = 1 and $||A|| ||A^{-1}|| \ge 1$ for all matrix norms induced by a vector norm.

Exercise 7.10 Let $x = (1, -3, 7)^T$. Compute $||x||_1, ||x||_2$, and $||x||_{\infty}$.

Exercise 7.11 Let $\bar{x} = (1.23, 0.37, -2.6)^T$ and assume that the elements \bar{x}_k are correctly rounded. Compute both the absolute and relative error measured in $\|\cdot\|_{\infty}$.

Exercise 7.12 Let

$$A = \left(\begin{array}{rrrr} 1.2 & 0.3 & -2.7 \\ 3.1 & -0.5 & 3.2 \\ 1.6 & -0.8 & -2.3 \end{array}\right).$$

Compute $||A||_{\infty}$.

Exercise 7.13 Formulate the following system of equations on matrix form Ax = b:

$$\begin{cases} 2x_1 - x_2 + 2x_3 = 1, \\ 3x_1 + x_2 - x_3 = 0, \\ -3x_1 + x_2 + 2x_3 = -3 \end{cases}$$

Exercise 7.14 Find a permutation matrix P such that

$$P\begin{pmatrix} x_1\\ x_2\\ x_3\\ x_4 \end{pmatrix} = \begin{pmatrix} x_1\\ x_4\\ x_2\\ x_3 \end{pmatrix}.$$

Exercise 7.15 Find a permutation matrix P such that

Exercise 7.16 Find a Gauss transformation M such that

$$M\left(\begin{array}{c}2\\3\\0.6\\-1.8\end{array}\right) = \left(\begin{array}{c}2\\3\\0\\0\end{array}\right).$$

Exercise 7.17 Suppose we have a linear system Ax = b where

$$A = \begin{pmatrix} 2 & 1 & -2 \\ -1 & 0 & 3 \\ 1 & 2 & -1 \end{pmatrix} \text{ and } b = \begin{pmatrix} 6 \\ 1 \\ -3 \end{pmatrix}.$$

During the first step of Gaussian elimination we multiply the system with a matrix M_1 such that the new system $M_1Ax = M_1b$ is

$$\left(\begin{array}{ccc} 2 & 1 & -2 \\ 0 & 0.5 & 2 \\ 0 & 1.5 & 0 \end{array}\right) x = \left(\begin{array}{c} 6 \\ 4 \\ -6 \end{array}\right).$$

Give the Gausstransformation M_1 .

Exercise 7.18 A Gauss transformation M_1 that eliminates the non-zeros from the first column of a matrix has the form

$$M_1 = \left(\begin{array}{rrrr} 1 & 0 & 0 \\ -m_{21} & 1 & 0 \\ -m_{31} & 0 & 1 \end{array}\right).$$

Show that the matrix can be written as $M_1 = I - me_1^T$. Also give the elements of the vector m in terms of the elements of the matrix $A = (a_{ij})$ and show that $M_1^{-1} = I + me_1^T$.

Exercise 7.19 The Gauss transformation M_2 that is used during the second step of Gaussian elimination has the structure

$$M_2 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & -m_{32} & 1 & 0 \\ 0 & -m_{42} & 0 & 1 \end{pmatrix}.$$

Show that M_2 can be written as $M_2 = I - me_2^T$ where $m = (0, 0, m_{32}, m_{42})$. Also show that $M_2^{-1} = I + me_2^T$.

Exercise 7.20 Let

$$A = \begin{pmatrix} 2 & 1 & -2 \\ -3 & 0.5 & -2 \\ 1 & 2.5 & 0 \end{pmatrix}.$$

Is pivoting required during the first step of Gaussian elimination? If so give the appropriate permutation matrix to use.

Exercise 7.21 After one step of Guassian elimination we have

$$\left(\begin{array}{rrrr} 2 & 1 & -2 \\ 0 & -0.5 & 1.6 \\ 0 & 1.7 & 0.3 \end{array}\right).$$

Give the permutation matrix P_2 and the Gausstransformation M_2 to use the next step of the Guassian elimination. Also carry out the step and given the resulting upper triangular matrix U.

Exercise 7.22 We need to solve a system of equations Ax = b, where

$$A = \begin{pmatrix} 1 & 1.2 & -2.3 \\ 2 & 0 & -1 \\ -1 & 1.5 & 2.1 \end{pmatrix} \text{ and } b = \begin{pmatrix} -0.5 \\ 1.2 \\ 1.27 \end{pmatrix}.$$

We compute the decomposition PA = LU where

$$L = \begin{pmatrix} 1 & 0 & 0 \\ -0.5 & 1 & 0 \\ 0.5 & 0.8 & 1 \end{pmatrix} \quad U = \begin{pmatrix} 2 & 0 & -1 \\ 0 & 1.5 & 1.6 \\ 0 & 0 & -3.08 \end{pmatrix} \quad P = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}.$$

Take advantage of the decomposition to compute the solution x.

Exercise 7.23 A computer program has computed the decomposition PA = LU and the output is

$$L = \begin{pmatrix} 1 & 0 & 0 \\ -0.7 & 1 & 0 \\ 0.3 & 1.8 & 1 \end{pmatrix} \quad U = \begin{pmatrix} 1.7 & -2.3 & -1.4 \\ 0 & 1.2 & -0.5 \\ 0 & 0 & 3.1 \end{pmatrix} \quad P = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}.$$

Determine if pivoting was used correctly during the computations.

Exercise 7.24 Let a matrix A and a vector \overline{b} be given as

$$A = \begin{pmatrix} 1 & -2 & 0 \\ 2 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix} \qquad \bar{b} = \begin{pmatrix} -1.03 \\ 1.34 \\ 0.78 \end{pmatrix}$$

where the elements of \bar{b} are correctly rounded. The system of equations Ax = b was solved using Gaussian elimination and the solution $x = (0.1867, -0.4217, 0.9667)^T$. Find an upper bound for the relative error in the solution x caused by the rounding errors in the right hand side.

Hint You may use that $||A^{-1}||_{\infty} = 1$.

Exercise 7.25 Let

$$A = \begin{pmatrix} 3.1 & -1.2 & 2.6\\ 1.5 & -0.7 & 3.6\\ -4.1 & 1.1 & 0.7 \end{pmatrix}.$$

If we solve the system $Ax = b_1$, where $b = (0.654, 0.765, -1.042)^T$ we obtain the solution $x = (1.2595, 3.5113, 0.3705)^T$. Give an upper bound for the change in the solution x if we instead use the approximate right hand side $\bar{b} = (0.657, 0.761, -1.039)^T$.

Hint Use that $||A^{-1}||_{\infty} = 17.9052.$

Exercise 7.26 Let PA = LU be the LU decomposition. Prove the formula

$$\det(A) = (-1)^k \prod_{i=1}^n u_{ii}.$$

What is k here?

Exercise 7.27 Let $r = b - A\hat{x}$ be the residual for an approximate solution to the linear system Ax = b. Prove the formula:

$$||x - \hat{x}|| \le ||A^{-1}|| ||r||.$$

Exercise 7.28 Let

$$A = \begin{pmatrix} 2 & 1 & -1 \\ 1 & 3 & -2 \\ 2 & 0 & 1 \end{pmatrix} \text{ och } = \begin{pmatrix} 1.12 \\ -1.07 \\ 2.34 \end{pmatrix}.$$

and the decomposition PA = LU, where

$$L = \begin{pmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ 1 & -0.4 & 1 \end{pmatrix} \quad U = \begin{pmatrix} 2 & 1 & -1 \\ 0 & 2.5 & -1.5 \\ 0 & 0 & 1.4 \end{pmatrix} \quad P = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

be given. Use the decomposition to compute the determinant det(A). You may use that det(AB) = det(A)det(B).

8 Least Squares Problems and Orthogonal Decompositions

Exercise 8.1 Suppose we have a series of measurements (x_i, y_i) , i = 1, 2, ..., m, and want to find a function of the type $y = f(x) = c_0 + c_1 x + c_2 \sin(x)$ that best fits the data in the least squares sense. Formulate the problem as an over determined linear system Ax = b.

Exercise 8.2 An harmonic wave is described by the amplitude and the phase shift, i.e. $u(t) = A \cdot \sin(t + \phi)$. Suppose we have a sequence of measurements $(t_i, u(t_i))$, i = 1, 2, ..., m, and want to find A and ϕ using the least squares method. First explain why the least squares method cannot be applied directly. Also rewrite the problem in such a way that the least squares method can be used.

Exercise 8.3 The time needed for a certain algorithm to complete can be described by the formula $t \approx Cn^p$, where *n* is the problem size, *p* is the computational complexity, and *C* is the average time needed for one arithmetic operation. The following table is available

n	100	200	300	400
t	1.213	2.370	3.619	4.875

Formulate the problem of estimating the parameters C and p from the given data as an over determined linear system Ax = b.

Exercise 8.4 Suppose Q is an orthogonal matrix. Show that $||Qx||_2 = ||x||_2$, for all vectors x, and thus $||Q||_2 = 1$.

Exercise 8.5 If A is both an ortgohonal matrix and an orthogonal projection. What can you conclude about A?

Exercise 8.6 Suppose $A \in \mathbb{R}^{m \times n}$, m > n, and that we have the reduced QR decomposition $A = Q_1 R$. Show how the decomposition can be used to find the vector x that minimize $||Ax - b||_2$.

Exercise 8.7 Suppose $A \in \mathbb{R}^{m \times n}$, m > n, and that we have the QR decomposition

$$A = Q \left(\begin{array}{c} R\\ 0 \end{array} \right) = Q_1 R$$

Show that the linear system Ax = b has an exact solution if $b = Q_1 Q_1^T b$.

Exercise 8.8 Consider the vector a as an $n \times 1$ matrix. Write out its reduced QR decomposition explicitly. Also write down a formula for the solution of the least squares problem $ax \approx b$, where b is a given $n \times 1$ vector.

Exercise 8.9 We are intressted in the least squares problem min $||Ax - b||_2$. Suppose $A = Q_1 R$ is the reduced QR decomposition. Use Q_1 to give a formula for a orthogonal projection P, such that Pb = r = b - Ax, where x is the least squares solution.

Exercise 8.10 Compute the reduced QR factorization of the matrix

$$A = \left(\begin{array}{cc} 0 & \sqrt{2} \\ -1 & 1 \\ 1 & 1 \end{array}\right).$$

Exercise 8.11 Let $W \in \mathbb{R}^{n \times n}$ be real, symmetric, positive definite, and let $\|\cdot\|_W$ be defined by,

$$||x||_W^2 = x^T W x.$$

Verify that $||x||_W = 0$ if and only if x = 0. Also derive the normal equations for the minimization problem,

$$\min_{x} \|Ax - b\|_{W}.$$

Hint Use the Cholesky factorization $W = R^T R$.

Exercise 8.12 Show that $||A||_2 = \sigma_1$ and if A^{-1} exists then $||A^{-1}||_2 = 1/\sigma_n$.

Exercise 8.13 Let A^T be an $m \times n$ matrix of rank $k < \min(m, n)$. Use the decomposition $A = U\Sigma V^T$ to give an orthogonal basis for null (A^T) .

Exercise 8.14 Show that if $A \in \mathbb{R}^{m \times n}$ has rank *n*, then $||A(A^T A)^{-1}A^T||_2 = 1$.

Exercise 8.15 Suppose the matrix $B \in \mathbb{R}^{m \times n}$ has full column rank. Use the decomposition $B = U\Sigma V^T$ to give a formula for the solution to the the problem

$$\min_{x} ||Bx||_2$$
, subject to $||x||_2 = 1$.

Exercise 8.16 Suppose $A \in \mathbb{R}^{m \times n}$, m > n, rank(A) = n, and that we have a factorization $A = U\Sigma V^T$. Clearly demonstrate how the matrices U and V provides basis vectors for the spaces Range(A) and null(A). What are the dimension of the range and null space respectively.

Exercise 8.17 Let $A \in \mathbb{R}^{m \times n}$, m > n, and $\operatorname{rank}(A) = n$. Demonstrate how the decomposition $A = U\Sigma V^T$ can be used for solving the least squares problem

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2$$

Give formulas for both the solution x and the redisual r = b - Ax.

Exercise 8.18 Let $A \in \mathbb{R}^{m \times n}$, m < n, and $\operatorname{rank}(A) = m$. Let $b \in \mathbb{R}^m$. Show that the formula

$$x = \sum_{i=1}^{m} \frac{u_i^T b}{\sigma_i} v_i$$

provides a solution to Ax = b. Is the solution unique?

Exercise 8.19 Suppose we want to find the solution to a linear system Ax = b, where rank(A) = k < n so that the solution x is not unique. Demonstrate how the solution x can be split into two parts,

 $x = x_1 + x_2,$ $x_1 \in \operatorname{null}(A)^{\perp},$ and, $x_2 \in \operatorname{null}(A),$

and how the SVD of A can be used to write expressions for the solution components x_1 and x_2 .

Exercise 8.20 Consider the Least Squares problem with linear constraints,

 $\min \|Ax - b\|_2, \quad \text{for all } x \in \mathbb{R}^n \text{ such that } Bx = 0,$

where A is $m \times n$, m > n, and B is $n \times n$.

- a) Suppose $\operatorname{rank}(B) = n$. What is the solution of the least squares problem?
- b) Suppose rank(B) = k < n. Show how the SVD can be used to derive a formula for the solution of the least squares problem.

Exercise 8.21 Let $A \in \mathbb{R}^{m \times n}$, where $m \gg n$, have full column rank. Use the decomposition $A = U\Sigma V^T$ to develop a criteria that ensures that the linear system Ax = b has a solution. Try and make the criteria as inexpensive as possible to check.

Exercise 8.22 Tikhonov regularization means replacing an ill-conditioned linear system Ax = b by the more stable problem,

$$\min_{x} \|Ax - b\|_{2}^{2} + \lambda^{2} \|x\|_{2}^{2},$$

where λ is the regularization parameter. Show that the normal equations of the above least squares problem are

$$(A^T A + \lambda^2 I)x = A^T b.$$

Also derive a formula for the singular values of the matrix $(A^T A + \lambda^2 I)$ and use the result to show that the normal equations are not ill-conditioned (provided λ is selected appropriately). Finally derive a formula for the solution x_{λ} .

1 Basic Concepts and Floating Point systems

Exercise 1.1 We first observe that $|\Delta a| \leq 0.5 \cdot 10^{-4}$ and $|\Delta b| \leq 0.4 \cdot 10^{-2} < 0.5 \cdot 10^{-2}$. Thus a has 4 correct decimals and b has 2. Further a has 3 significant digits and b has 4.

Exercise 1.2 We rewrite \bar{c}_0 as 2997924.58 Thus if the approximate value is roundned correctly it has two correct decimals. With 7 digits in the integer part the total number of significant digits is this 9.

Exercise 1.3 In order to obtain 5 significant digits we need 3 correct decimals. Thus $\bar{a} = 22.735$ and $|\Delta a| \leq 0.5 \cdot 10^{-3}$.

Exercise 1.4 The absolute error in the approximation is $|\pi - \bar{\pi}| \leq 9.3 \cdot 10^{-5} = 0.093 \cdot 10^{-3} < 0.5 \cdot 10^{-3}$. Thus the approximation has 3 correct decimals and 4 significant digits.

Exercise 1.5 The rounded value 1.61803 has five correct decimals because $|\phi - 1.61803| = 0.00000398875... = 0.398875... \times 10^{-5} \le 0.5 \times 10^{-5}$. Further, the the rounded value 1.6180 has five significant digits.

Exercise 1.6 Note that $|\Delta a| \leq 4.7 \cdot 10^{-3} < 0.5 \cdot 10^{-2}$. Hence \bar{a} has 2 correct decimals. Add the two digits before the decimal point and we have 4 significant digits.

Exercise 1.7 The rounding of π is an error in used data R_X , which is propagated to the result. The approximation of the exponential function is a truncation error R_T .

Exercise 1.8 The unit round-off for the number system is $\mu = 0.5 \cdot 10^{-3}$. Thus $|\Delta x| \le \mu |x| \le 0.5 \cdot 10^{-3} 102.232 < 0.052$.

Exercise 1.9 First write as a normalized number $x = 1.683556541 \cdot 10^2$. Then round to 5 correct digits in the fractional part and obtain $\bar{x} = 1.68356 \cdot 10^2$.

Exercise 1.10 The numbers are stored as follows



Exercise 1.11 We rewrite the numbers in normalized form and obtain $x = 1.175614 \cdot 10^2$ and $y = 1.678214 \cdot 10^{-2}$. If we round the fractional parts to 5 digits we find $x_r = 1.17561 \cdot 10^2$ and $y_r = 1.67821 \cdot 10^{-2}$.

Exercise 1.12 We say that m is the mantissa, β is the base (or radix) of the number system, and e is the exponent. To fully define a floating point number system we would also need to define the precision t to prescribe the number of digits in the mantissa $m = \pm d_0.d_1d_2...d_t$. Also, the floating point system needs lower L and upper U bounds for the value of the exponent, i.e. $L \leq e \leq U$

Exercise 1.13 First do the calculations exactly

 $a \cdot b = 4.4625 \cdot 10^8$ and $a/b = 1.6190476 \dots \cdot 10^3$.

Then round the results to two digits in the fractional part to obtain

$$fl[a \cdot b] = 4.46 \cdot 10^8$$
 and $fl[a/b] = 1.62 \dots \cdot 10^3$.

In both cases a bound for the relative error is the unit round off for the floating point system, i.e. $\mu = 0.5 \cdot 10^{-2}$.

Exercise 1.14 For the first expression we rewrite

$$\sqrt{1+x} - 1 = \frac{(\sqrt{1+x}-1)(\sqrt{1+x}+1)}{\sqrt{1+x}+1} = \frac{x}{\sqrt{1+x}+1}$$

For the second

$$\frac{1}{1-x} - \frac{1}{1+x} = \frac{(1+x) - (1-x)}{(1+x)(1-x)} = \frac{2x}{1-x^2}$$

which avoids the cancellation if x is small. For the last expression we use

$$1 - \cos^2(x) = \sin^2(x).$$

Exercise 1.15 Since there is no error in the number 10 and absolute errors are added during a minus operation the error in the computed value \bar{x} is also at most $0.5 \cdot 10^{-5}$. A smaller error is achived by avoiding the cancellation, i.e. use $x = 1/(\sqrt{101} + 10)$

2 Error Analysis

Exercise 2.1 First compute the approximate value $\overline{f} = 1.8$ ($|R_B| = 0$). Since x, y and z are correctly rounded we have error bounds $|\Delta x|, |\Delta y|, |\Delta z| \leq 0.5 \cdot 10^{-2}$. The error propagation formula gives us

$$|\Delta f| \lesssim |\frac{\partial f}{\partial x}||\Delta x| + |\frac{\partial f}{\partial y}||\Delta y| + |\frac{\partial f}{\partial z}||\Delta z| = |\frac{1}{z}||\Delta x| + |\frac{-1}{z}||\Delta y| + |-\frac{x-y}{z^2}||\Delta z| \le 0.5 \cdot 10^{-2}.$$

Thus $f = 1.80 \pm 0.5 \cdot 10^{-2}$.

Exercise 2.2 The error propagation formula gives

$$|\Delta f(2)| \le |\frac{\partial f(2)}{\partial b}| |\Delta b| = |2e^{2b}| |\Delta b| = |2e^{2 \cdot 1.70}| 0.01 < 0.6.$$

The approximate value is $f(2) = 29.9641... \approx 30.0$, with $|R_B| \leq 0.5 \cdot 10^{-1}$. Add both errors to obtain $f(2) = 30.0 \pm 0.7$.

Exercise 2.3 The approximate area is $\bar{A} = \bar{\pi}\bar{r}^2 = 1773.77713920 \approx 1774$, $|R_B| \leq 0.3$. Note that the error in $\bar{\pi}$ is at most $0.5 \cdot 10^{-3}$. The error propagation formula gives

$$|\Delta A| \lesssim |\frac{\partial A}{\partial \pi}| |\Delta \pi| + |\frac{\partial A}{\partial r}| |\Delta r| = |r^2| |\Delta \pi| + |2\pi r| |\Delta r| \approx 3.2684 < 3.3.$$

The total error is $|R_{TOT}| \le 3.3 + 0.3 < 4$. Thus $A = 1774 \pm 4$.

Exercise 2.4 The approximate value is $\bar{y} = 1.54e^{3.17} = 36.7$, $|R_B| \le 0.5 \cdot 10^{-1}$. The error propagation formula gives

$$|\Delta y| \lesssim |\frac{\partial y}{\partial a}||\Delta a| + |\frac{\partial y}{\partial b}||\Delta b| = |\mathbf{e}^b||\Delta a| + |a\mathbf{e}^b||\Delta b| < 2.55.$$

The total error is $|R_{TOT}| \le 2.55 + 0.5 \cdot 10^{-1} < 2.6$. Thus $y = 36.7 \pm 2.6$.

In hindsight it would probably have been better to round to 37 and use $|R_B| = |37 - 36.6635...| < 0.34$ to obtain the answer $y = 37 \pm 3$. Either works fine.

Exercise 2.5 Rewrite the expression to read

$$f(a,b) = \frac{ab}{a+b}.$$

The approximate value is $\bar{f} = 18.9$, $|R_B| \leq 0.05$, and the error propagation formula gives

$$|\Delta f| \lesssim |\frac{\partial f}{\partial a}||\Delta a| + |\frac{\partial f}{\partial b}||\Delta b| = |\frac{b^2}{(a+b)^2}||\Delta a| + |\frac{a^2}{(a+b)^2}||\Delta b| \le 0.57.$$

The total error is $|R_{TOT}| \le 0.05 + 0.57 < 0.7$. Thus $f = 18.9 \pm 0.7$.

Exercise 2.6 We first determine the computational order as

$$f(x) = \sqrt{1+x} - \sqrt{1-x} = \sqrt{a} - \sqrt{b} = c - d = e.$$

The relative errors in the intermediate results, e.g. $|\Delta a|/|a|$, are bounded by μ . The error propagation formula gives

$$|\Delta f| \lesssim |\frac{1}{2\sqrt{a}}||\Delta a| + |\frac{1}{2\sqrt{b}}||\Delta b| + |\Delta c| + |\Delta d| + |\Delta e|.$$

In order to simplify the result we use $a \approx b \approx c \approx d \approx 1$ for small x. Also

$$f(x) = \frac{(\sqrt{1+x} - \sqrt{1-x})(\sqrt{1+x} + \sqrt{1-x})}{\sqrt{1+x} - \sqrt{1-x}} = \frac{2x}{\sqrt{1+x} - \sqrt{1-x}} \approx x,$$

for small x. We obtain

$$\Delta f| \lesssim \mu(\frac{1}{2} + \frac{1}{2} + 1 + 1 + |x|) \approx 3\mu.$$

Since $f(x) \approx x$ for small x the bound for the relative error is $|\Delta f|/|f| \leq 3|x|^{-1}\mu$.

Exercise 2.7 In the first case we have the computational order

$$f_1(x) = \frac{1 - \cos(x)}{\sin(x)} = \frac{1 - c}{s} = \frac{d}{s} = e.$$

The error propagation formula gives us

$$\begin{split} |\Delta f| \lesssim |\frac{\partial f}{\partial c}||\Delta c| + |\frac{\partial f}{\partial s}||\Delta s| + |\frac{\partial f}{\partial d}||\Delta d| + |\frac{\partial f}{\partial e}||\Delta e| = |\frac{1}{s}||\Delta c| + |\frac{1-c}{s^2}||\Delta s| + |\frac{1}{s}||\Delta d| + |\Delta e| \leq \\ \mu(|\frac{c}{s}| + |\frac{1-c}{s}| + |\frac{d}{s}| + |e|) \approx \frac{\mu}{x} \end{split}$$

where we have used $c \approx 1$, $s \approx x$ and $d/s = e = f \approx x/2$. Similarly for the second expression we use the computational order

$$f_2(x) = \frac{\sin(x)}{1 + \cos(x)} = \frac{s}{1 + c} = \frac{s}{d} = e.$$

The error propagation formula gives us

$$|\Delta f| \lesssim |\frac{s}{(1+c)^2}||\Delta c| + |\frac{1}{1+c}||\Delta s| + |\frac{s}{d^2}||\Delta d| + |\Delta e| \le \mu(|\frac{cs}{(1+c)^2}| + |\frac{s}{1+c}| + |\frac{s}{d}| + |e|) \approx 1.75x\mu.$$

These are the absolute errors. Insert $x = 1.111 \cdot 10^{-8}$ in and $|\Delta f_1| \le 10^{-8}$ which is on the same order of magnitude as the actual difference between f_1 and f_2 .

Exercise 2.8 The computational order is

$$f(x) = e^x - 3x = a - 3x = a - b = c$$

The error propagation formula gives us

$$|\Delta f| \lesssim |\frac{\partial f}{\partial a}||\Delta a| + |\frac{\partial f}{\partial b}||\Delta b| + |\frac{\partial f}{\partial c}||\Delta c| = |1||\Delta a| + |1||\Delta b| + |1||\Delta c| \lesssim |\Delta f| \leq |\Delta f| \leq |\Delta f| + |\Delta f$$

$$\mu(|a| + |b| + |c|) \approx \mu(|1| + |3x| + |1|) \approx 2\mu,$$

where we have used $e^x \approx 1$, $f(x) = c \approx 1$ since x is small. There is no cancellation present in these calculations. Everything turns out fine and both the absolute and relative errors are bounded by 2μ (since the function value $f(x) \approx 1$).

Exercise 2.9 The computational order and the intermediate results are

$$f = \frac{x - \sin(x)}{x^3} = \frac{x - a}{b} = \frac{c}{b} = d.$$

where the relative error in all the intermediate results are bounded by μ . The error propagation formula gives

$$\begin{split} |\Delta f| \lesssim |\frac{\partial f}{\partial a}||\Delta a| + |\frac{\partial f}{\partial b}||\Delta b| + |\frac{\partial f}{\partial c}||\Delta c| + |\frac{\partial f}{\partial d}||\Delta d| = \\ |-\frac{1}{b}||\Delta a| + |-\frac{c}{b^2}||\Delta b| + |\frac{1}{b}||\Delta c| + |1||\Delta d| \le \mu(|\frac{a}{b}| + |\frac{c}{b}||\frac{c}{b}| + |d|) \approx \mu(|\frac{1}{x^2}| + \frac{3}{6}). \end{split}$$

If we insert $x = 10^{-7}$ and $\mu \approx 1.1 \cdot 10^{-16}$ we obtain $||\Delta f| \approx 0.0110$. Thus the actual error $5.4 \cdot 10^{-3}$ is within the error bound.

Exercise 2.10 The computational order is

$$f(x) = 1 - 2x\cos(x) = 1 - 2xa + 1 - b = c.$$

The error propagation formula gives us

$$\begin{split} |\Delta f| \lesssim |\frac{\partial f}{\partial a}||\Delta a| + |\frac{\partial f}{\partial b}||\Delta b| + |\frac{\partial f}{\partial c}||\Delta c| &= |2x||\Delta a| + |1||\Delta b| + |1||\Delta c| \lesssim \\ \mu(|2xa| + |b| + |c|) \approx \mu(|2x| + |2x| + 1) \approx \mu, \end{split}$$

where we have used $\cos(x) \approx 1$, $f(x) = c \approx 1$ and that x is small. There is no cancellation present in these calculations. Everything turns out fine and both the absolute and relative errors are bounded by μ (since the function value $f(x) \approx 1$).

Exercise 2.11 We do a Taylor series expansion of $f(\bar{x})$ around x to obtain

$$f(\bar{x}) = f(x + \Delta x) = f(x) + f'(x)\Delta x + f''(\eta)\frac{(\Delta x)^2}{2}.$$

where $\eta \in (x, \bar{x})$ and f'(x) = 0. If Δx is small then $\eta \approx \bar{x}$ and we obtain

$$|\Delta f| \lesssim |f''(\bar{x})| \frac{|\Delta x|^2}{2}.$$

3 Non-linear equations

Exercise 3.1 Since the function is continuous and f(0) = -1 and f(1) = 1.3849 there has to be a root in the interval [0, 1].

Exercise 3.2 $\phi'(x) = \frac{1}{2}e^x$ and $|\phi'(0.3)| \approx 0.68 < 1$.

Exercise 3.3 With $\bar{x} = 1.7$ we get $f(\bar{x}) = -0.387$. Also $f'(x) = 3x^2 + 1$ so $f'(\xi) \approx f'(1.7) = 9.67$. Thus

$$|1.7 - x^*| \le \frac{|f(1.7)|}{|f'(1.7)|} < \frac{0.39}{9.6} < 0.0406 < 0.5 \cdot 10^{-1}.$$

Exercise 3.4 Since $f(\bar{x})$ is evaluated as zero we need to include the computational errors and we actually have $\bar{f}(\bar{x}) = 0$ and $|\bar{f}(\bar{x}) - f(\bar{x})| \le 10^{-8}$. The error estimate is

$$|\bar{x} - x^*| \le \frac{|f(\bar{x})|}{|f'(\xi)|} \le \frac{|\bar{f}(\bar{x})| + |\bar{f}(\bar{x}) - f(\bar{x})|}{|f'(\xi)|} \le \frac{0 + 10^{-8}}{1.7} < 6 \cdot 10^{-9}.$$

Exercise 3.5 Let $f(x) = 5(1 - e^{-x}) - x$ and find a root of the equation f(x) = 0. We use the Newton-Raphson method. If $x_0 = 5$ we get $x_2 = \bar{x} = 4.96511$. The error estimate is

$$|\bar{x} - x^*| \le \frac{|f(\bar{x})|}{|f'(\bar{x})|} \le \frac{4.1 \cdot 10^{-6}}{0.96} \le 0.5 \cdot 10^{-5}.$$

Thus $x^* = 4.96511 \pm 0.5 \cdot 10^{-5}$ has five correct decimals.

Exercise 3.6 The rate of convergence is determined by the derivative $|\phi'(x^*)|$. For the two methods, and $x^* \approx 1.9$, we get $\phi'_1(1.9) \approx 0.647$ and $\phi'_2(1.9) \approx 0.177$. Thus the second method has the fastest convergence. Also note that a fixed point satisfies $x^* = \frac{x^*}{2} + \sin(x^*)$ which can be written as $x^* = \sin(x^*)$. Thus a fixed point satisfies the original equation and we really have convergence to a root.

Exercise 3.7 Let x^* be the fixed point. Then

$$|x_n - x^*| = |\phi(x_{n-1}) - \phi(x^*)| \le |\phi'(\xi_n)| |x_{n-1} - x^*| \le C |x_{n-1} - x^*|,$$

where the constant C satisfies $|\phi'(\xi_n)| \leq C$. In this case $\phi(x) = \frac{1}{2}\cos(x)$ and $\phi'(x) = -\frac{1}{2}\sin(x)$. Thus $|\phi'(x)| \leq \frac{1}{2}$. Thus we can pick $C = \frac{1}{2} < 1$ and we have proved that the error $|x_n - x^*|$ is reduced by a factor of 2 in each step. Thus the method is convergent regardless of x_0 .

In order to obtain the error estimate, for $\bar{x} = 0.4501832$, we comute the derivative $f'(x) = -2\sin(x) - 4$ and use the estimate = |-4.87026...| > 4.8 = M. Thus

$$|\bar{x} - x^*| \le \frac{|f(\bar{x})|}{|f'(\bar{x})|} \le \frac{2.01 \cdot 10^{-6}}{4.87} < 0.42 \cdot 10^{-6}.$$

We conclude that $x^* = 0.4501832 \pm 0.5 \cdot 10^{-6}$.

Exercise 3.8 First the error in the approximation $\bar{x} = 1.05$ is estimated by

$$|\bar{x} - x^*| \le \frac{|f(\bar{x})|}{|f'(\bar{x}|)} \approx \frac{|f(1.05)|}{|f'(1.05)|} < \frac{1.9 \cdot 10^{-4}}{2.0} < 10^{-4},$$

where $f'(x) = 1 + 3e^{-x}$.

Now let x^* be a fixed point to the iterations. In the case (i) the fixed point obviously satisfies the equation f(x) = 0. In the case (ii) we get $3x^* = 2x^* + 3e^{-x^*}$ or $x^* = 3e^{-x^*}$ or $f(x^*) = 0$. A similar resoning holds in the case (iv). However, a fixed point to (iii) does not satisfy the equation.

Now we consider the convergence speed of the iterations (i), (ii) and (iv). In all cases we compute the derivative of the iteration function and obtain

$$\phi_1(x) = 3e^{-x}$$
, so $\phi'_1(1.05) = 1.05$ and divergence,
 $\phi_2(x) = (2x + 3e^{-x})/3$, so $\phi'_2(1.05) = -0.32$ and convergence,

and

$$\phi_4(x) = (x + 3e^{-x})/2$$
, so $\phi'_4(1.05) = -0.025$, and again convergence,

Thus the method *(iv)* should converge the fastest to x^* .

Finally the error in step k is $|x_k - x^*| \approx |\phi'_4(1.05)|^k |1.05 - x^*| \approx (0.025)^k 10^{-4}$. Testing shows that k = 4 gives an error of $3.9 \cdot 10^{-11}$ and thus 4 iterations is enough.

Exercise 3.9 If we apply Newton-Raphson to the equation $x^2 = 0$ we get

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} = x_k - \frac{x_k^2}{2x_k} = \frac{1}{2}x_k,$$

which means that if $x_0 = 1$ then $x_k = 2^{-k} \to 0 = x^*$ as $k \to \infty$. The rate of convergence is verified to be linear by the observation that $e_k = |x_k - x^*| = |\frac{1}{2}x_{k-1} - 0| = \frac{1}{2}|x_{k-1} - 0| = \frac{1}{2}e_{k-1}$.

Exercise 3.10 Newton-Raphsons method is defined by the iteration function

$$\phi(x) = x - \frac{f(x)}{f'(x)}$$
, and $\phi'(x) = -\frac{f(x)f''(x)}{(f'(x))^2}$.

Since x^* is a single root, i.e. $f'(x^*) \neq 0$, we see that $\phi'(x^*) = 0$. A Taylor series expansion shows that

$$\phi(x_k) = \phi(x^*) + \phi'(x^*)(x_k - x^*) + \frac{\phi''(\xi)}{2}(x_k - x^*)^2, \xi \in (x_k, x^*).$$

Since $\phi(x_k) = x_{k+1}$, $\phi(x^*) = x^*$ and $\phi'(x^*) = 0$ we obtain

$$x_{k+1} - x^* = \frac{\phi''(\xi)}{2}(x_k - x^*)^2,$$

which shows that the convergence is quadratic.

Exercise 3.11 We use a Taylor series expansion

$$f(\bar{x}) = f(x^*) + f'(x^*)(\bar{x} - x^*) + \frac{f''(\xi)}{2}(\bar{x} - x^*)^2, \xi \in (\bar{x}, x^*),$$

and since x^* is a double root we have $f(x^*) = f'(x^*) = 0$ and thus

$$f(\bar{x}) = \frac{f''(\xi)}{2}(\bar{x} - x^*)^2.$$

Taking absolute values and approximating $\xi \approx \bar{x}$ gives the desired estimate. In the practical case when $\bar{x} = 1.99789$ we find that $f(\bar{x}) = -6.6875 \cdot 10^{-6}$ and $f''(\bar{x}) = -3.0127$. Thus

$$|\bar{x} - x^*| \le \sqrt{\frac{2|f(\bar{x})|}{|f''(\bar{x})|}} \le \sqrt{\frac{2 \cdot 6.7 \cdot 10^{-6}}{3.0}} < 2.2 \cdot 10^{-3}.$$

Note that $x^* = 2$ so the actual error is of the same magnitude as the error estimate in this case.

Exercise 3.12 First we use the error estimate

$$|\bar{x} - x^*| \lesssim \frac{|f(\bar{x})|}{|f'(\bar{x})|} \approx \frac{1.2971 \cdot 10^{-4}}{3.0433} < 4.3 \cdot 10^{-5}.$$

Secondly, we define the order of convergence as the largest integer p such that

$$\lim_{k \to \infty} \frac{|x_k - x^*|}{|x_{k-1} - x^*|^p} = C < \infty.$$

Since x^* is unknown we cannot directly apply the definition. The simplest solution is to assume that the iteration x_4 has a much smaller error than the other iterations x_1, x_2, x_3 . Thus we approximate $x^* = 0.5177574$ and compute the errors $|x_0 - x^*| \approx 4.8 \cdot 10^{-1}$, $|x_1 - x^*| \approx 1.4 \cdot 10^{-1}$, $|x_2 - x^*| \approx 1.4 \cdot 10^{-2}$, and $|x_3 - x^*| \approx 1.5 \cdot 10^{-4}$. Since $(|x_1 - x^*|)^2 \approx (1.4 \cdot 10^{-1})^2 \approx 2 \cdot 10^{-2} \approx |x_2 - x^*|$ and $(|x_2 - x^*|)^2 \approx (1.4 \cdot 10^{-2})^2 \approx |x_3 - x^*|$ we conclude that the table shows that p = 2 for Newton-Raphsons method.

Exercise 3.13 The order of convergence is the largest integer p such that

$$\lim_{k \to \infty} \frac{|x_k - x^*|}{|x_{k-1} - x^*|^p} = C < \infty.$$

This means that $|x_k - x^*| \approx C |x_{k-1} - x^*|^p$, for a constant *C*. In the table we see that the error in step k+1 is always about half of the error at step k. This fits very nicely with p = 1 and C = 0.5. Thus we have *linear convergence* for this specific polynomial. It is known that Newton-Raphsons method has order of convergence p = 2 for single roots and p = 1 for double roots and this $x^* = 2$ has to be a *double root*.

Exercise 3.14 First if we apply the Newton-Raphson method to $f(x) = x^2 - a = 0$ we obtain

$$x_{k+1} = x_k - \frac{x_k^2 - a}{2x_k} = \frac{1}{2}(x_k - a/x_k) = \phi(x_k).$$

Since a is a normalized floating point number we can write $a = (-1)^2 (1.f)_2 2^k$. We only need to consider positive numbers and if k is even then $\sqrt{a} = \sqrt{(1.f)_2 2^{k/2}}$ and if k is odd we instead have $\sqrt{a} = \sqrt{(1.f)_2 2^{1/2} 2^{(k-1)/2}}$. Thus in worst case we need to compute the square root of a number $1 \le (1.f)_2 2^1 < 4$. So it is enough to consider the case $1 \le a < 4$ and $1 \le x < 2$.

For the convergence analysis we compute $\phi''(x) = -\frac{a}{x^3}$ and therefore $\phi''(\sqrt{a}) = 1/\sqrt{a} \le 1$. Thus

$$|x_k - \sqrt{a}| \lesssim \frac{1}{2} |x_{k-1} - \sqrt{a}|^2$$

If $x_0 = 1.5$ the maximum error is $|x_0 - \sqrt{a}| \le 0.5$. We get $|x_1 - \sqrt{a}| \le 0.125$, $|x_2 - \sqrt{a}| \le 0.078$, $|x_3 - \sqrt{a}| \le 3.1 \cdot 10^{-5}$, $|x_4 - \sqrt{a}| \le 4.7 \cdot 10^{-10}$, and $|x_5 - \sqrt{a}| \le 1.1 \cdot 10^{-19}$. We see that 5 iterations are needed if $x_0 = 1.5$.

In order to decrease the number of iterations by one we need an initial guess with an error $|x_0 - \sqrt{a}| < 0.125$. Then the new x_0 is the same as x_1 above. This means dividing the interval [1, 4] into n intervals $[a_0, a_1)$, $[a_1, a_2), \ldots, [a_{n-2}, a_{n-1})$. For a certain starting value a we indentify the index k such that $a_k \leq a < a_{k+1}$. The initial guess x_0 is then given by the middle point $(\sqrt{a_k} + \sqrt{a_{k+1}})/2$ and the initial error is $(\sqrt{a_{+1}} - \sqrt{a_k})/2$. If n = 4 then the largest initial error is given by $(\sqrt{1 + 3/n} - \sqrt{1})/2 = 0.1614$. If we continue as above we get $|x_4 - \sqrt{a}| \leq 6.5 \cdot 10^{-18}$ which is below machine precision.

4 Polynomial and Spline Interpolation

Exercise 4.1 A polynomial, $p_n(x)$, of degree n, can be written $p_n(x) = c_0 + c_1x + \ldots + c_nx^n$. Thus there are n + 1 parameters to determine and the polynomial can satisfy exactly n + 1 interpolation conditions.

Exercise 4.2 Use the Ansatz $p(x) = c_0 + c_1(x - 0.9) + c_2(x - 0.9)(x - 1.1)$, where the last term is used to estimate the truncation error. The interpolation conditions give

$$p(0.9) = c_0 = f(0.9) = 0.4710$$
, and $p(1.1) = c_0 + c_1(1.1 - 0.9) = f(1.1) = 0.2452$ so $c_1 = -1.1290$.

Thus the linear polynomial is p(x) = 0.4710 - 1.1290(x - 0.9). The truncation error is obtained using

$$p(1.2) = c_0 + c_1(1.2 - 0.9) + c_2(1.2 - 0.9)(1.2 - 1.1) = f(1.2) = 0.2385$$
 so $c_2 = 3.5400$.

The truncation error is $R_T \approx 3.45(x-0.9)(x-1.1)$. Insert x = 1.03 to find $f(1.03) \approx p(1.03) = 0.3242$, $|R_B| \leq 0.5 \cdot 10^{-4}$, and $|R_T(1.03)| \leq 0.033$. The errors in the table also gives an error $|R_{XF}| \leq 0.5 \cdot 10^{-4}$ and we obtain $f(1.03) = 0.3242 \pm 0.034$. Its resonable to round off a bit more to get $f(1.03) = 0.324 \pm 0.04$.

Exercise 4.3 We require that the total error is $R_{TOT} \leq 10^{-5}$. If the values in the table is stored with 6 correct digits then the resulting error is $R_{XF} \leq 0.5 \cdot 10^{-6}$. Thus the truncation error can be at most $R_T \leq 9.5 \cdot 10^{-6}$. For linear interpolation the truncation error is given by

$$R_T \le \frac{h^2}{8} \max_{x_i \le \xi \le x_{i+1}} |f''(\xi)|,$$

and since $f(x) = \log(x)$ we find that $|f''(x)| = |-x^{-2}| \le 1$ since $1 \le x \le 4$. Thus $R_T \le 9.5 \cdot 10^{-6}$ if $h^2 \le 8 \cdot 9.5 \cdot 10^{-6}$ which gives $h \le 0.0087$. Since the interval length is $x_n - x_1 = 4 - 1 = 3$ the required table size is $n = 3/h \approx 3/0.0087 \approx 344.82 < 345$. Thus we need at least 345 function values in our table.

Exercise 4.4 The error estimate for linear interpolation is

$$f(x) - p(x) = \frac{1}{2}f''(\eta(x))(x - x_1)(x - x_2),$$

where, for our case, $x_1 = 0$, $x_2 = 1$, and $0 < \eta(x) < 1$. Since $f(x) = \sin(x)$ we obtain $|f''(\eta)| = |-\sin(\eta)| \le 1$. We also see that $|x(x-1)| \le \frac{1}{4}$ since the maximum occurs for $x = \frac{1}{2}$. Thus

$$|f(x) - p(x)| = \frac{1}{2}|f''(\eta(x))(x - 0)(x - 1)| \le \frac{1}{2}1\frac{1}{4} = \frac{1}{8}.$$

Exercise 4.5 Use the Ansatz $p(x) = c_0 + c_1(x - 0.0) + c_2(x - 0.0)(x - 0.5)$, where the last term is used to estimate the truncation error. The interpolation conditions p(0.0) = 1.80 and p(0.5) = 2.80 gives us the linear interpolating polynomial $p_1(x) = 1.8 + 2x$. The final interpolation condition p(1.0) = 4.10 gives us truncation error $R_T \approx 0.6(x - 0.0)(x - 0.5)$. Insert x = 0.4 to obtain $f(1.03) \approx p_1(0.4) = 2.6$, where $|R_B| = 0$, and $|R_T(0.4)| \le 0.024$. The errors in the table also gives an error $|R_{XF}| \le 0.5 \cdot 10^{-2}$ and we obtain $f(0.4) = 2.6 \pm 0.03$.

Exercise 4.6 Use the Ansatz $p(x) = c_0 + c_1(x - 0.7) + c_2(x - 0.7)(x - 0.8) + c_3(x - 0.7)(x - 0.8)(x - 0.6)$, where the last term is used to estimate the truncation error. The interpolation conditions p(0.7) = 1.29, p(0.8) = 1.32 and p(0.6) = 1.23 gives the coefficients of the quadratic polynomial $c = (1.29, 0.3, -1.5)^T$. The truncation error is obtained from the interpolation condition p(0.9) = 1.07. We find that $c_3 \approx -41.7$. We find that $p_2(0.74) = 1.306$, with $|R_B| \le 0.5 \cdot 10^{-3}$ (not asked for in the exercise), and with the truncation error $|R_T| \le -0.014$.

Exercise 4.7 The basis function satisfies $\ell_2(x_2) = 1$ and $\ell_2(x_i) = 0$, $i \neq 2$. Thus

$$\ell_2(x) = \frac{(x-x_1)(x-x_3)(x-x_4)}{(x_2-x_1)(x_2-x_3)(x_2-x_4)}$$

The degree of $\ell_2(x)$ is n = 3.

Exercise 4.8 The polynomial is

$$p(x) = 1.3 \frac{(x-2)(x-3)}{(1-2)(1-3)} + 0.6 \frac{(x-1)(x-3)}{(2-1)(2-3)} + 1.9 \frac{(x-1)(x-2)}{(3-1)(3-2)}.$$

There is no reason to simplify the expression further.

Exercise 4.9 First $p(0) = c_0 = 0$ and $p(1) = c_0 + c_1 + c_2 + c_3 = 0$ gives two equations. Then $p'(x) = c_1 + 2c_2x + 3c_3x^2$ so we also obtain $p'(0) = c_1 = 1$ and $p'(1) = c_1 + 2c_2 + 3c_3 = 1$. Thus the system of equations is

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 2 & 3 \end{pmatrix} \begin{pmatrix} c_0 \\ c_1 \\ c_2 \\ c_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \end{pmatrix}.$$

We can solve the linear system by noting that $c_0 = 0$ and $c_1 = 1$. Then we are left with two equations for c_2 and c_3 . The solution is $p(x) = x - 3x^2 + 2x^3$.

Exercise 4.10 The conditions for s(x) to be a cubic spline are (i) on each sub interval $[x_i, x_{i+1}]$ the spline s(x) should be given by a cubic polynomial, and (ii) s(x), s'(x) and s''(x) should be continuous on the whole interval $[x_1, x_n]$. Also the (iii) the interpolation conditions $s(x_i) = f(x_i)$ needs to be satisfied. The given information is not sufficient since we also need two end point conditions for the spline to be unique.

Exercise 4.11 The function s(x) is a cubic spline since s(x), s'(x) and s''(x) are continuous at x = 1.

Exercise 4.12 The conditions that has to be satisfied are s(1) = a + 1 = b + c, s'(1) = a = 3b + 2c and s''(1) = 0 = 6b + 2c. The solution is a = -3, b = 1 and c = -3.

Exercise 4.13 Since f(x) is a cubic polynomial, and correct end point conditions are to be used, then s(x) = f(x). This means the spline s(x) is the same cubic polynomial in each of the sub intervals.

Exercise 4.14 We have more than enough information to determine the parameters. The most efficient way is to first compute $s_2(1) = 2$ and determine a using $s_1(1) = 1.6 + a = 2$, i.e. a = 0.4. Now compute $s'_1(1) = 2.5$ and use $s'_2(1) = b = 2.5$. This leaves c which can be computed using $s_2(2) = 4.9 + c = 6.7$ or c = 1.8.

Exercise 4.15 Check the continuity requirements by $s_1(1) = 0.6 = s_2(1)$, $s'_1(1) = -0.6 = s'_2(1)$ and finally $s''_1(1) = -0.8 = s''_2(1)$. Thus s(x) is a cubic spline. We also compute $s''_1(0) = -0.8$ and conclude that s(x) is not a natural cubic spline.

Exercise 4.16 For s(x) to be a cubic spline we require that $s_{k-1}(x_k) = s_k(x_k)$, $s'_{k-1}(x_k) = s'_k(x_k)$ and $s''_{k-1}(x_k) = s''_k(x_k)$, for k = 1, 2, ..., N - 1. In addition we need the interpolation conditions $s_k(x_k) = f(x_k)$ and $s_k(x_{k+1}) = f(x_{k+1})$, for k = 0, 1, ..., N - 1. A sketch of the case N = 3 is given below



Exercise 4.17 Since the error satisfies $E(N) \approx Ch^p = C'(1/N)^p$, where $C' = C(b-a)^p$ is a constant. Thus $E(N)/E(2N) = 2^p$ and if we insert the values from the table we get

$$\frac{E(5)}{E(10)} = \frac{0.0435}{0.00269} \approx 16.2 \text{ and } \frac{E(10)}{E(20)} = \frac{0.00269}{0.000172} \approx 15.6$$

In both cases the quitient is sufficiently close to $2^4 = 16$ to conclude that p = 4.

Exercise 4.18 The function s(x) is a linear spline if it is given by a first degree polynomial on each subinterval. In our case the intervals are 1 < x < 3 and 3 < x < 4. Thus we are seeking two linear polynomials $s_1(x)$ and $s_2(x)$. We can write $s_1(x) = 1.56 + c_1(x-1)$ and use $s_1(3) = 2.31$ to obtain $c_1 = 0.375$. Similarly we obtain $s_2(x) = 2.31 - 0.34(x-3)$. The linear spline is thus given by

$$s(x) = \begin{cases} 1.56 + 0.375(x-1), & 1 \le x < 3, \\ 2.31 - 0.34(x-3), & 3 \le x < 4. \end{cases}$$

Exercise 4.19 For a) differentiate the expression for p(t) to obtain the tangent

$$p'(t) = -2(1-t)P_1 + 2(1-2t)P_2 + 2tP_3$$
, and $p'(0) = 2(P_2 - P_1)$.

For **b**) we first note that P_4 and P_5 should have the same *x*-coordinate so $P_4 = (6, , \alpha)^T$. The tangent direction at P_3 should be parallel to $P_3 - P_2 = (3, 5)^T - (2, 6)^T = (1, -1)^T$. Compute $P_4 - P_3 = (6, \alpha)^T - (3, 5)^T = (3, \alpha - 5)^T = 3(1, -1)^T$ if $\alpha = 2$. Thus we have to use $P_4 = (6, 2)^T$.

Exercise 4.20 The sketch is



The convex hull is the area enclosed by the dashed lines. Important features of the Beziér curve is that since both P_1/P_2 and P_3/P_4 have the same x-coordinate the tangent direction of the curve is vertical at both the starting and ending points.

Exercise 4.21 For a) we differentiate

$$p'(t) = -3(1-t)^2 P_1 + 3(-2(1-t)t + (1-t)^2) P_2 + 3(2(1-t)t - t^2) P_3 + 3t^2 P_4, \text{ and } p'(0) = 3(P_2 - P_1).$$

For **b**) the definition of the *convex hull* is the set of all convex linear combinations $\alpha_1 P_1 + \alpha_2 P_2 + \alpha_3 P_3 + \alpha_4 P_4$, where $\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 = 1$ and $\alpha_1, \alpha_2, \alpha_3, \alpha_4 \ge 0$. To show that the Beziér curve is located within the convex hull we note that the weights in the expression for p(t) are calculated from the identity

$$1 = 1^{3} = (1 - t + t)^{3} = (1 - t)^{3} + 3(1 - t)^{2}t + 3(1 - t)t^{2} + t^{3},$$

and that all terms are positive for $0 \le t \le 1$. For c) we compute p(1/2) by inserting t = 1/2 into the expression (note that t = 0.5 means also 1 - t = 0.5)

$$p(1/2) = \left(\frac{1}{2}\right)^3 \left(\left(\begin{array}{c} 0\\0 \end{array}\right) + 3 \left(\begin{array}{c} 1\\3 \end{array}\right) + 3 \left(\begin{array}{c} 4\\2 \end{array}\right) + \left(\begin{array}{c} 5\\1 \end{array}\right) \right) = \left(\begin{array}{c} 2.5\\2.0 \end{array}\right)$$

The sketch is



Exercise 4.22 For the slope to be -2 the tangent vector at P_4 should be in the direction $(1, -2)^T$. Thus we can pick

$$P_3 = \begin{pmatrix} 2\\3 \end{pmatrix} - \alpha \begin{pmatrix} 1\\-2 \end{pmatrix}$$
 and $P_5 = \begin{pmatrix} 2\\3 \end{pmatrix} + \alpha \begin{pmatrix} 1\\-2 \end{pmatrix}$,

where α is a positive number. There is no unique solution to this problem.

5 Integration, Differentiation and Extrapolation

Exercise 5.1 Taylors formula gives

$$f(x+h) = f(x) + f'(x)h + \frac{1}{2}f''(x)h^2 + \frac{1}{6}f^{(3)}(x)h^3 + \dots$$

Thus

$$\frac{1}{h}(f(x+h) - f(x)) = f'(x) + \frac{1}{2}f''(x)h + \frac{1}{6}f^{(3)}(x)h^2 + \dots = f'(x) + Ch + \mathcal{O}(h^2), \quad C = \frac{1}{2}f''(x).$$

Exercise 5.2 First we compute two different approximations using the Trapezoidal method. We have

$$T(0.5) = 0.5\left(\frac{1.000}{2} + 1.006 + \frac{0.670}{2}\right) = 0.9205,$$

and

$$T(0.25) = 0.25 \left(\frac{1.000}{2} + 1.015 + 1.006 + 0.882 + \frac{0.670}{2}\right) = 0.9345.$$

Since the Trapezoidal method has a truncation error of the form $R_T = ch^2$ we can estimate the truncation error by

$$R_T(0.25) = \frac{T(0.5) - T(0.25)}{3} = -4.667 \cdot 10^{-3}.$$

Also since the table is given with absolute error $|\Delta f| < 0.5 \cdot 10^{-3}$ we get an error in the integral that can be estimated $|R_{XF}| \le (1-0)0.5 \cdot 10^{-3} = 0.5 \cdot 10^{-3}$. If we round the result to 3 correct digits we get $I = 0.934 \pm (0.5 \cdot 10^{-3} + 4.7 \cdot 10^{-3} + 0.5 \cdot 10^{-3}) = 0.934 \pm 6 \cdot 10^{-3}$.

Exercise 5.3 Since $|\Delta f(x)|/|f(x)| \leq \varepsilon$ we obtain the error estimate

$$|R_{XF}| = |\int_a^b f(x)dx - \int_a^b \bar{f}(x)dx| \le \int_a^b |f(x) - \bar{f}(x)|dx \le \int_a^b \varepsilon |f(x)|dx \le \left(\int_a^b |f(x)|dx\right)\varepsilon.$$

This means that the bound is good for the case when f(x) is either positive or negative. If f(x) changes sign in the interval (a, b) then it may happen that

$$\left|\int_{a}^{b} f(x)dx\right| << \int_{a}^{b} |f(x)|dx|$$

and the result may have a large relative error.

Exercise 5.4 We denote the error by $\epsilon_h \approx Ch^p$. Then

$$\frac{\epsilon_{h_1}}{\epsilon_{h_2}} \approx \frac{Ch_1^p}{Ch_2^p} = (\frac{h_1}{h_2})^p$$

Insert numbers from the table we obtain

$$2^p = \left(\frac{0.2}{0.1}\right)^p \approx \frac{\epsilon_{0.2}}{\epsilon_{0.1}} = \frac{0.342}{0.0861} \approx 3.97 \text{ and } 2^p \approx \frac{\epsilon_{0.1}}{\epsilon_{0.05}} = \frac{0.0861}{0.0209} \approx 4.11.$$

We see that $2^p = 4$ which means p = 2.

Exercise 5.5 We see that

$$F_{1}(h) + \frac{F_{1}(h) - F_{1}(qh)}{q^{p_{1}} - 1} = a + bh^{p_{1}} + ch^{p_{2}} + \frac{a + bh^{p_{1}} + ch^{p_{2}} - (a + bq^{p_{1}}h^{p_{1}} + cq^{p_{2}}h^{p_{2}})}{q^{p_{1}} - 1} = a + bh^{p_{1}} + ch^{p_{2}} + \frac{b(1 - q^{p_{1}})h^{p_{1}} + c(1 - q^{p_{2}})h^{p_{2}}}{q^{p_{1}} - 1} = a + c\left(1 - \frac{1 - q^{p_{2}}}{1 - q^{p_{1}}}\right)h^{p_{2}} = a + \mathcal{O}(h^{p_{2}})$$

Exercise 5.6 Let h = 0.05. From the table we see that

$$T(h) = I + R_T \approx I + Ch^4$$
, and $T(2h) \approx I + C(2h)^4 = I + 16Ch^4 \approx I + 16R_T$

where I is the exact value of the integral and R_T is the truncation error for h = 0.05. Thus

$$R_T \approx Ch^4 \approx \frac{T(2h) - T(h)}{15} = 7.83 \cdot 10^{-6}.$$

Exercise 5.7 With $T(h) = T_0 + Ch^p$ we form the quotient

$$\frac{T(h) - T(h/2)}{T(h/2) - T(h/4)} = \frac{h^p(1 - 2^{-p})}{(h/2)^p(1 - 2^{-p})} = 2^p.$$

Insert the values for h = 0.4 and h = 0.2 to obtain $2^p \approx 4.11$ and $2^p \approx 3.73$. We conclude that p = 2.

Exercise 5.8 First we compute

$$T(1/2) = 0.5\left(\frac{1.5000}{2} + 1.0713 + \frac{0.9073}{2}\right) = 1.137475,$$

and

$$T(1/4) = 0.25 \left(\frac{1.5000}{2} + 1.2412 + 1.0713 + 0.9663 + \frac{0.9073}{2}\right) = 1.1206125.$$

The truncation error in the result for h = 1/4 is estimated by $|R_T| \leq |T(1/2) - T(1/4)|/3 < 5.7 \cdot 10^{-3}$ since the truncation error for the Trapezoidal method has the form $R_T \approx ch^2$. Since the function values are correctly rounded we also have $|R_{XF}| \leq (b-a)\varepsilon = (1-0)0.5 \cdot 10^{-4}$. This means that the integral can be estimated by $I = 1.121 \pm (0.5 \cdot 10^{-4} + 5.7 \cdot 10^{-3} + 0.5 \cdot 10^{-3}) = 1.121 \pm 7 \cdot 10^{-3}$.

Exercise 5.9 We find that

$$\frac{R_T(h)}{R_T(h/2)} = 2^p.$$

with h = 0.1 we get $2^p \approx 1.98$ and with h = 0.05 we get $2^p \approx 1.98$. Thus p = 1 and for h = 0.025we get $R_T = -0.0099 = C(0.025)^1$ which means $C \approx -0.4$. The truncation error is thus given by $R_T \approx -0.4h^1$ and $R_T < 10^{-8}$ if $h < 10^{-8}/0.4 = 2.5 \cdot 10^{-8}$.

Exercise 5.10 No since the Trapezoidal method needs, at least a continuous second derivative f''(x) on the interval [1,4] to have a truncation error $R_T \approx ch^2$. In this case the first derivative is

$$f'(x) = -2x\sin(x^2)\sqrt{4-x} - \cos(x^2)(4-x)^{-1/2}$$

which tends to $-\infty$ as $x \to 4$. This means that f(x) does not have the nessecary smoothness for the Trapezoidal method to work well.

Exercise 5.11 For a) we compute

$$\frac{S(4h) - S(2h)}{S(2h) - S(h)} \approx \frac{c4^p h^p - c2^p h^p}{c2^p h^p - h^p} = 2^p.$$

From the table, with h = 1/8, we obtain

$$\frac{S(1/2) - S(1/4)}{S(1/4) - S(1/8)} \approx \frac{2.06246 - 2.03208}{2.03208 - 2.03021} = 16.25 \approx 2^4.$$

We thus conclude that p = 4. We can then use

$$S(1/4) - S(1/8) = 2.03208 - 2.03021 \approx c(2^4 - 1)(1/8)^4,$$

to obtain c = 0.51. For **b**) we note that $|R_T| \approx 0.51h^4$ is smaller than 10^{-6} if $h \leq (10^{-6}/0.51)^{1/4} \approx 0.03742$. Thus a step size h < 0.037 is required.

6 Differential Equations

Exercise 6.1 With $y(0) = y_0 = 1$, and h = 0.5 we obtain

$$y(0.5) \approx y_1 = y_0 + hf(t_0, y_0) = y_0 + h(y_0^2 - t_0) = 1 + 0.5(1^2 - 0) = 1.5.$$

Also, since the truncation error for Euler's method is $R_T \approx ch$ we find $y(0.5; 5h) - y(0.5; h) \approx y(0.5) + c5h - y(0.5) - ch = 4ch$. Thus, with h = 0.1 we get $R_T \approx c5h \approx 5(y(0.5; 5h) - y(0.5; h))/4 \approx 0.021$.

Exercise 6.2 Using Euler's method we find that, using h = 0.2, we have $y(0.2; h = 0.2) = y_1 = y_0 + hf(t_0, y_0) = 1.6$. If we instead use h = 0.1 we get

$$y(0.1) \approx y_1 = y_0 + hf(t_0, y_0) = 2 - 0.1 \cdot (0^2 - 2) = 1.8$$
 and $y(0.2) \approx y_2 = y_1 + hf(t_1, y_1) \approx 1.621$

Considering that the truncation error for Euler is $R_T = Ch$ we also find $y(0.2; 2h) - y(0.2; h) \approx y(0.2) + c2h - y(0.2) - ch = ch$. Thus, with h = 0.1 we get $R_T \approx ch \approx (y(0.2; 2h) - y(0.2; h)) \approx 0.021$. For the complete error estimate we also note that the value y(0.2; h = 0.1) = 1.621 is rounded to 4 correct decimals. This $|R_B| \leq 0.5 \cdot 10^{-4}$ and $y(0.2) = 1.621 \pm 0.022$.

Exercise 6.3 Euler's method applied to the problem gives the $y_{k+1} = (1 - 100h)y_k + \sin^2(t_k)$. For stability we only look at the homogeneous part, i.e. $y_{k+1} = (1 - 100h)y_k$. Thus the requirement is $-1 \le 1 - 100h \le 1$, or, since h is positive, $h \le 2/100 = 0.02$.

Exercise 6.4 By applying the Euler method we obtain the difference formula $y_{k+1} = (1-100h)y_k + h \sin^2(t_k)$. Thus the homogeneous part is stable if $|1 - 100h| \le 1$, or $h \le 2/100 = 0.02$.

Exercise 6.5 The method is *implicit* since the right hand side has the term $f(t_{k+1}, y_{k+1})$ and y_{k+1} is the unknown we seek to compute. Thus we must solve an equation in each step. For the stability we apply the method to the test problem $y' = \lambda y$, y(0) = 1, and obtain $y_{k+1} = y_k + \frac{h\lambda}{2}(y_k + y_{k+1})$ or

$$y_{k+1} = \left(\frac{1 + \frac{h\lambda}{2}}{1 - \frac{h\lambda}{2}}\right) y_k$$

Thus the Trapezoidal method is stable if $|1 + \frac{h\lambda}{2}| \le |1 - \frac{h\lambda}{2}|$. This means that $h\lambda/2$ is closer to -1 than to 1. Since h is positive this means that the method is stable for all $\lambda < 0$. So if the solution $y(t) = e^{\lambda t}$, to the test problem, should be decreasing, i.e. for $\lambda < 0$, then the numerical solution will also be decreasing.

Exercise 6.6 Since $f(t, y) = \lambda y$ we obtain

$$k_1 = hf(t_k, y_k) = h\lambda y_k$$
, and $k_2 = hf(t_k, y_k + k_1) = h\lambda(y_k + k_1) = h\lambda(y_k + h\lambda y_k) = (h\lambda + (h\lambda)^2)y_k$.

The next iterate y_{k+1} is thus given by

$$y_{k+1} = y_k + \frac{1}{2}(k_1 + k_2) = y_k + \frac{1}{2}(h\lambda + h\lambda + (h\lambda)^2)y_k = (1 + h\lambda + (h\lambda)^2/2)y_k.$$

Thus the method is stable if $|1 + h\lambda + (h\lambda)^2/2| \le 1$.

Exercise 6.7 We define v = y'. Then v' = y'' and

$$\begin{pmatrix} y \\ v \end{pmatrix}' = \begin{pmatrix} v \\ ty - 3v - 1 \end{pmatrix}, \quad \begin{pmatrix} y(0) \\ v(0) \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

Exercise 6.8 Introduce v = y'. Then $v' = y'' = 2y(1 + y^2)$. Thus the system is

$$\begin{pmatrix} y \\ v \end{pmatrix}' = \begin{pmatrix} v \\ 2y(1+y^2) \end{pmatrix}, \quad \begin{pmatrix} y(0) \\ v(0) \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}.$$

If we apply the Euler method we get, with h = 0.2,

$$\begin{pmatrix} y_1 \\ v_1 \end{pmatrix} = \begin{pmatrix} y_0 \\ v_0 \end{pmatrix} + h \begin{pmatrix} v_0 \\ 2y_0(1+y_0^2) \end{pmatrix} = \begin{pmatrix} 1.4 \\ 2.8 \end{pmatrix}.$$

Similarly $(y_2, v_2)^T = (1.96, 4.576)^T$. Thus we obtain the approximation $y(0.4) \approx y_2 = 1.96$. Exercise 6.9 Introduce v = y' to obtain

$$v' + e(y^2 - a)v + y = 0$$
, $y(0) = a$, $v(0) = b$.

The system is thus

$$\left(\begin{array}{c} y\\ v\end{array}\right)' = \left(\begin{array}{c} v\\ -e(y^2-a)v-y\end{array}\right), \quad \left(\begin{array}{c} y(0)\\ v(0)\end{array}\right) = \left(\begin{array}{c} a\\ b\end{array}\right).$$

7 Basic matrix operations and linear systems of equations

Exercise 7.1 We evaluate the expression using the following operations

$$z = (A + I)Bx + y = (A + I)x_1 + y = Ax_1 + x_1 + y = x_2 + x_1 + y = x_3 + y = x_4$$

Computing the matrix vector product $x_1 = Bx$ requires mn multiplications and additions each, i.e. a total of 2mn operations. The product $x_2 = Ax_1$ requires $2m^2$ operations. The remaining two vector additions require m additions (as $y, x_1 \in \mathbb{R}^m$). So the operation count is m(2m + 2n + 2).

Exercise 7.2 Aim to keep intermediate results small. Multiplication by an inverse is dealt with by solving the corresponding linear system, i.e. compute $z = A^{-1}x$ by solving Az = x. The order of computation is

$$z_1 = Ab$$
, $Cz_2 = b$, $z_3 = z_1 + z_2$, $z_4 = Az_3$, $z_5 = 2z_4 + z_3$, and finally $Bx = z_5$.

All intermediate results are vectors.

Exercise 7.3 An upper triangular system Rx = y can be solved using,

$$x_i = \left(b_i - \sum_{j=i+1}^n r_{ij} x_j\right) / r_{ii}, \quad i = n, n - 1, \dots, 1.$$

Thus in step i exactly (i - 1) multiplications and additions are needed. Also exactly one division. Thus the total amount of work is

$$\sum_{i=1}^{n} (i-1) \approx \frac{n^2}{2}$$

multiplications and additions or n^2 arithmetic operations.

Exercise 7.4 The *LU*-decomposition and two triangular systems require $2n^3/3 + 2n^2$ arithmetic operations. For large *n* we have $t(n) \approx c \cdot n^3$, where *c* is the average time for one operation, and thus

$$\frac{t(2n)}{t(n)} \approx \frac{c2^3n^3}{cn^3} = 8$$

If we compute the same quotients using the table we obtain

$$\begin{array}{c|cccc} n & 1024 & 2048 & 4096 \\ \hline t(2n)/t(n) & 7.4166 & 7.1835 & 7.8295 \end{array}$$

The results look promising and fits quite well with the idea that the computational complexity is $\mathcal{O}(n^3)$.

Exercise 7.5 The unit circles are



Exercise 7.6 A matrix norm is *induced* if its definition is based on a vector norm, i.e.

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

For such norms we have

$$\|AB\| = \max_{x \neq 0} \frac{\|ABx\|}{\|x\|} = \max_{x \neq 0} \frac{\|ABx\|}{\|Bx\|} \frac{\|Bx\|}{\|x\|} \le \left(\max_{x \neq 0} \frac{\|ABx\|}{\|Bx\|}\right) \left(\max_{x \neq 0} \frac{\|Bx\|}{\|x\|}\right) \le \max_{y \neq 0} \frac{\|Ay\|}{\|y\|} \|B\| \le \|A\| \|B\|.$$

Exercise 7.7 Demonstrate the first inequality by

$$||x||_{\infty}^{2} = \max_{1 \le i \le n} |x_{i}|^{2} \le \sum_{i=1}^{n} |x_{i}|^{2} = ||x||_{2}^{2}.$$

Also, since $|x_i| \leq ||x||_{\infty}$, we have

$$||x||_2^2 = \sum_{i=1}^n |x_i|^2 \le \sum_{i=1}^n ||x||_\infty^2 = n ||x||_\infty.$$

Exercise 7.8 Recall the definition

$$||uv^T||_2 = \max_{x \in \mathbb{R}^n} \frac{||uv^Tx||_2}{||x||_2} = \max_{x \in \mathbb{R}^n} \frac{|v^Tx|||u||_2}{||x||_2}$$

The Cauchy-Schwarz inequality is $|v^T x| \leq ||v||_2 ||x||_2$ with equality for x = v. So

$$||uv^{T}||_{2} = \frac{|v^{T}v|||u||_{2}}{||v||_{2}} = ||v||_{2}||u||_{2}.$$

Exercise 7.9 First from the definition of the matrix norm, and since Ix = x we have

$$||I|| = \max_{x \neq 0} \frac{||Ix||}{||x||} = \max_{x \neq 0} \frac{||x||}{||x||} = 1$$
, so $1 = ||I|| = ||AA^{-1}|| \le ||A|| ||A^{-1}||$.

Exercise 7.10 If $x = (1, -3, 7)^T$ then $||x||_{\infty} = 7$, $||x||_2 = \sqrt{59}$ and $||x||_1 = 11$.

Exercise 7.11 If $\bar{x} = (1.23, 0.37, -2.6)^T$ is correctly rounded then the error vector satisfies $|\delta x| \le (0.005, 0.005, 0.05)^T$. Thus $||x - \bar{x}||_{\infty} \le 0.5 \cdot 10^{-1}$ is the absolute error and $||x - \bar{x}||_{\infty}/||x||_{\infty} \le 0.05/2.6 < 0.02$ is the relative error.

Exercise 7.12 From the second row we get the largest sum and $||A||_{\infty} = 3.1 + 0.5 + 3.2 = 6.8$.

Exercise 7.13 The linear system is Ax = b where

$$A = \begin{pmatrix} 2 & -1 & 2 \\ 3 & 1 & -1 \\ -3 & 1 & 2 \end{pmatrix} \text{ and } b = \begin{pmatrix} 1 \\ 0 \\ -3 \end{pmatrix}.$$

Exercise 7.14 The rows should appear in the order 1 - 4 - 2 - 3 and thus

$$P = \left(\begin{array}{rrrrr} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{array}\right).$$

Exercise 7.15 The rows should appear in the order 4 - 3 - 1 - 2 which gives

$$P = \left(\begin{array}{rrrr} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{array}\right).$$

Exercise 7.16 The multipliers are $m_3 = 0.6/3 = 0.2$ and $m_4 = -1.8/3 = -0.6$. Therefore the Gauss transformation is

$$M = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0.2 & 1 & 0 \\ 0 & -0.6 & 0 & 1 \end{pmatrix}.$$

Exercise 7.17 The multipliers are $m_{21} = -1/2 = -0.5$ and $m_{31} = 1/2 = 0.5$. Therefore the Gauss transformation is

$$M_1 = \left(\begin{array}{rrrr} 1 & 0 & 0\\ 0.5 & 1 & 0\\ -0.5 & 0 & 1 \end{array}\right).$$

Exercise 7.18 Set $m = (0, m_{21}, m_{31})^T$ and find that

$$M_1 = I - me_1^T = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} - \begin{pmatrix} 0 \\ m_{21} \\ m_{31} \end{pmatrix} (1, 0, 0)$$

has the desired structure. To show that $M_1^{-1} = I + m e_1^T$ we do

$$(I - me_1^T)(I + me_1^T) = I - me_1^T + me_1^T - me_1^T me_1^T = I$$

where $e_1^T m = 0$ since *m* has a zero in the first position. The elements of *m* are chosen so that $0 = a_{i1} - m_{i1}a_{11}$ and thus $m_{i1} = a_{i1}/a_{11}$, i = 2, 3.

Exercise 7.19 Set $m = (0, 0, m_{32}, m_{42})^T$ and find that

$$M_{2} = I - me_{2}^{T} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} - \begin{pmatrix} 0 \\ 0 \\ m_{32} \\ m_{42} \end{pmatrix} (0, 1, 0, 0)$$

has the desired structure. To show that $M_2^{-1} = I + m e_2^T$ we do

$$(I - me_2^T)(I + me_2^T) = I - me_2^T + me_2^T - me_2^T me_2^T = I$$

where $e_2^T m = 0$ since m has zeros in the first two positions.

Exercise 7.20 Pivoting is required as $|a_{21}| > |a_{11}|$. The permutation

$$P_{12}\left(\begin{array}{rrr} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{array}\right)$$

should be used to switch rows one and two.

Exercise 7.21 We need to switch rows two and three to obtain

$$\left(\begin{array}{rrrr} 2 & 1 & -2 \\ 0 & 1.7 & 0.3 \\ 0 & -0.5 & 1.6 \end{array}\right).$$

and then use the multiplier $m_{32} = -0.5/1.7 = -0.2941$. The permutation matrix and Gausstransformation are

$$P_{23} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \text{ and } M_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0.2941 & 1 \end{pmatrix}.$$

The new element on position (3,3) will be $1.6 - (-0.2941) \cdot 0.3 = 1.6882$. The upper triangular matrix is

$$U = \left(\begin{array}{rrr} 2 & 1 & -2 \\ 0 & 1.7 & 0.3 \\ 0 & 0 & 1.6882 \end{array}\right)$$

Exercise 7.22 In order to take advantage of the decomposition we rewrite Ax = b as PAx = Pb and L(Ux) = Pb. By introducing the intermediate variable y = Ux we obtain two triangular systems Ly = Pb and Ux = y. First solve

$$Ly = Pb \text{ or } \begin{pmatrix} 1 & 0 & 0 \\ -0.5 & 1 & 0 \\ 0.5 & 0.8 & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} -0.5 \\ 1.2 \\ 1.27 \end{pmatrix} = \begin{pmatrix} 1.2 \\ 1.27 \\ -0.5 \end{pmatrix}$$

to obtain $y = (1.2, 1.87, -2.60)^T$. Next we solve Ux = y by backwards substitution and find that $x = (1.02, 0.34, 0.84)^T$.

Exercise 7.23 Since the multiplier $m_{32} = 1.8 > 1$ pivoting cannot have been used correctly.

Exercise 7.24 Since the elements of \bar{b} are correctly rounded to two correct decimals we have $\|\delta b\|_{\infty} \leq 0.5 \cdot 10^{-2}$. Thus the relative error in the solution is bounded by

$$\frac{\|\delta x\|_{\infty}}{\|x\|_{\infty}} \le \|A\|_{\infty} \|A^{-1}\|_{\infty} \frac{\|\delta b\|_{\infty}}{\|b\|_{\infty}} \le 3 \cdot 1 \cdot \frac{0.5 \cdot 10^{-2}}{1.34} < 0.012.$$

Exercise 7.25 The error in the right hand side is $\delta b = (0.003, 0.004, -0.003)^T$. We find that the error in the result is bounded by

$$\frac{\|\delta x\|_{\infty}}{\|x\|_{\infty}} \le \|A\|_{\infty} \|A^{-1}\|_{\infty} \frac{\|\delta b\|_{\infty}}{\|b\|_{\infty}} \le 6.9 \cdot 17.9052 \cdot \frac{0.004}{1.039} < 0.48$$

The absolute error is $\|\delta x\|_{\infty} \leq 0.48 \cdot 3.5113 < 1.7$.

Exercise 7.26 Use the LU decomposition of A to obtain

$$A = P^T L U$$
, so det $(A) = det(P^T) det(L) det(U)$.

Here both L and U are triangular so the determinant is the product of the diagonal elements. Also P is a permutation matrix. If we exchange two rows in a matrix then the determinant changes sign. So k is the number of row exchanges that actually occured during the Guassian elimination when computing the LU decomposition.

Exercise 7.27 Let $r = b - A\hat{x}$. Then $A^{-1}r = A^{-1}b - A^{-1}A\hat{x} = x - \hat{x}$ and we obtain $||x - \hat{x}|| \le ||A^{-1}|| ||r||$.

Exercise 7.28 Since P = I we have det(P) = 1 and we obtain

$$\det(A) = \det(P^T)\det(L)\det(U) = 1 \cdot 1 \cdot (2 \cdot 2.5 \cdot 1.4) = 7.$$

8 Least Squares Problems and Orthogonal Decompositions

Exercise 8.1 The model is $y_i = c_0 + c_1 x_i + c_2 \sin(x_i)$ and each measurement results in a row of a linear system Ax = b. We obtain

$$\begin{pmatrix} 1 & x_1 & \sin(x_1) \\ 1 & x_2 & \sin(x_2) \\ \vdots & \vdots & \vdots \\ 1 & x_m & \sin(x_m) \end{pmatrix} \begin{pmatrix} c_0 \\ c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}.$$

Exercise 8.2 The least squares method can not be used since the model is not linear. We rewrite the model as $u(t) = A\cos(\phi)\sin(t) + A\sin(\phi)\cos(t)$. The resulting over determined linear system is

$$\begin{pmatrix} \sin(t_1) & \cos(t_1) \\ \sin(t_2) & \cos(t_2) \\ \vdots & \vdots \\ \sin(t_m) & \cos(t_m) \end{pmatrix} \begin{pmatrix} A\cos(\phi) \\ A\sin(\phi) \end{pmatrix} = \begin{pmatrix} u(t_1) \\ u(t_2) \\ \vdots \\ u(t_m) \end{pmatrix}$$

where the first have to find a least squares solution x and then find the parameters A and ϕ by, e.g., $\phi = \operatorname{atan}(x_2/x_1)$.

Exercise 8.3 We rewrite the model as $\log(t_i) = \log(C) + p \log(n_i)$, and obtain the linear system

$$\begin{pmatrix} 1 & \log(100) \\ 1 & \log(200) \\ 1 & \log(300) \\ 1 & \log(400) \end{pmatrix} \begin{pmatrix} \log(C) \\ p \end{pmatrix} = \begin{pmatrix} \log(1.213) \\ \log(2.370) \\ \log(3.619) \\ \log(4.875) \end{pmatrix}$$

Exercise 8.4 Since Q is orthogonal $Q^T Q = I$. So

$$||Qx||_2^2 = (Qx)^T (Qx) = x^T Q^T Qx = x^T x = ||x||_2^2$$

This means that $||Q_2||_2 = \max ||Qx||_2 / ||x||_2 = \max 1 = 1.$

Exercise 8.5 First Range $(A) = \mathbb{R}^n$ since A is orthogonal and thus has linearly independent columns. So A is an orthogonal projection on the whole of \mathbb{R}^n . So Ax = x for every $x \in \mathbb{R}^n$ so A = I is the identity matrix.

Exercise 8.6 Let

$$A = Q \left(\begin{array}{c} R \\ 0 \end{array} \right) = Q_1 R$$

where $Q = (Q_1, Q_2)$. Since Q is orthogonal we find that

$$\|Ax - b\|_{2}^{2} = \|Q^{T}(Ax - b)\|_{2}^{2} = \|\begin{pmatrix} R\\ 0 \end{pmatrix} x - \begin{pmatrix} Q_{1}^{T}b\\ Q_{2}^{T}b \end{pmatrix}\|_{2}^{2} = \|Rx - Q_{1}^{T}b\|_{2}^{2} + \|Q_{2}^{T}b\|_{2}^{2}.$$

The minimum is achived for $x = R^{-1}Q_1^T b$. Thus only the reduced QR decomposition is needed.

Exercise 8.7 If $A = Q_1 R$ then range $(A) = \text{span}(Q_1)$. So $P = Q_1 Q_1^T$ is an orthogonal projection onto the range space of A. Thus if $b = Q_1 Q_1^T b$ then b belongs to the space range(A) which by definition means that Ax = b has a solution.

Exercise 8.8 The vector a can be seen as a matrix in $\mathbb{R}^{n \times 1}$. This means that

$$a = (a/||a||_2)||a||_2 = Q_1 R$$

where $Q_1 \in \mathbb{R}^{n \times 1}$ and $R \in \mathbb{R}^{1 \times 1}$. The formula for the least squares solution can be written using the normal equations $a^T a x = a^T b$ or $x = (a^T b)/(a^T a)$. This is the same as $x = R^{-1}Q_1^T b$ with the decomposition above.

Exercise 8.9 Since $x = R^{-1}Q_1^T b$ is the least squares solution we can compute the residual as $r = b - Ax = b - Q_1RR^{-1}Q_1^T b = (I - Q_1Q_1^T)b = Pb$. We can also recall that $Ax = Q_1Q_1^T b$ from the geometrical interpretation of the least squares problem.

Exercise 8.10 The matrix has orthogonal columns, i.e. if $A = (a_1, a_2)$ then $a_1^T a_2 = 0$. Thus the QR decomposition is

$$A = (a_1/\|a_1\|_2, \|a_2/\|a_2\|_2) \begin{pmatrix} \|a_1\|_2 & 0\\ 0 & \|a_2\|_2 \end{pmatrix} = Q_1 R.$$

The numbers are not very important.

Exercise 8.11 Let $W = R^T R$ be the Cholesky decomposition and rewrite

$$\|x\|_W^2 = x^T W x = x^T R^T R x = (Rx)^T (Rx) = \|Rx\|_2^2$$

Since R is non-singular Rx = 0 if and only if x = 0. In order to find the normal equations we use

$$||Ax - b||_W = ||R(Ax - b)||_2.$$

The normal equations are now $(RA)^T(RA)x = (RA)^T(Rb)$ or $A^TWAx = A^TWb$.

Exercise 8.12 If $A = U\Sigma V^T$ where U, V are orthogonal and $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n)$ then $||A||_2 = ||U\Sigma V^T||_2 = ||\Sigma||_2$. The norm of a diagonal matrix can be computed by

$$\|\Sigma\|_{2} = \max_{y \in \mathbb{R}^{n}} \frac{\|\Sigma y\|_{2}}{\|y\|_{2}} = \max_{y \in \mathbb{R}^{n}} \sqrt{\frac{\sum \sigma_{i}^{2} y_{i}^{2}}{\sum y_{i}^{2}}} \le \sigma_{1} \max_{y \in \mathbb{R}^{n}} \sqrt{\frac{\sum y_{i}^{2}}{\sum y_{i}^{2}}} = \sigma_{1}$$

with equality for $y = e_1$. Thus $||A||_2 = \sigma_1$. If A^{-1} exists then $A^{-1} = V\Sigma^{-1}U^T$ and $||A^{-1}||_2 = ||\Sigma^{-1}||_2$. Since the diagonal elements of Σ^{-1} are $1/\sigma_i$ the largest diagonal element is $1/\sigma_n$ and $||A^{-1}||_2 = 1/\sigma_n$.

Exercise 8.13 Let $A^T = V\Sigma^T U^T$. If $y \in \text{span}(u_{k+1}, \ldots, u_m)$ then $u_i^T y = 0$ for $i = 1, \ldots, k$. This is the null space of A^T .

Exercise 8.14 First compute $(A^T A)^{-1} = (V \Sigma^T U^T U \Sigma V^T)^{-1} = V(\Sigma^T \Sigma)^{-1} V^T$. Here $\Sigma^T \Sigma = \text{diag}(\sigma_i^2) \in \mathbb{R}^{n \times n}$. Thus $A(A^T A)^{-1} A^T = U \Sigma V^T V(\Sigma^T \Sigma)^{-1} V^T V \Sigma^T U^T = U \Sigma (\Sigma^T \Sigma)^{-1} \Sigma^T U^T$. Since U is orthogonal $||A(A^T A)^{-1} A||_2 = ||\Sigma(\Sigma^T \Sigma)^{-1} \Sigma^T||_2$. Evaluate the product of the diagonal matrices to obtain

$$\Sigma(\Sigma^T \Sigma)^{-1} \Sigma^T = \begin{pmatrix} I & 0 \\ 0 & 0 \end{pmatrix} \in \mathbb{R}^{m \times m}, \quad I \in \mathbb{R}^{n \times n}$$

The norm is the largest diagonal entry, i.e. 1.

Exercise 8.15 Let $B = U\Sigma V^T$. Since $V = (v_1, \ldots, v_n)$ provides a basis for \mathbb{R}^n any x can be written

$$x = \sum_{i=1}^{n} c_i v_i \Longrightarrow By = \sum_{i=1}^{n} c_i \sigma_i u_i.$$

If $||x||_2 = 1$ then $\sum c_i^2 = 1$. So

$$||Bx||_{2}^{2} = \sum_{i=1}^{n} \sigma_{i}^{2} c_{i}^{2} \ge \sigma_{n} \sum_{i=1}^{n} c_{i}^{2} = \sigma_{n}^{2},$$

with equality for $c = e_n$. So the minimum is σ_n and it is obtained for $x = \pm v_n$.

Exercise 8.16 The decomposition $A = U\Sigma V^T$ can be written

$$A = \sum_{i=1}^{n} \sigma_i u_i v_i^T,$$

where $\sigma_n > 0$ as rank(A) = n. This means that $Av_i = \sigma_i u_i \neq 0$ for i = 1, ..., n. So the null space is only the trivial one null $(A) = \{0\}$ with dimension 0. Similarly, if y belongs to the range then there is an x such that y = Ax, or

$$y = Ax = \sum_{i=1}^{n} \sigma_i(v_i^T x) u_i,$$

so the y is a linear combination of $\{u_1, \ldots, u_n\}$. Thus range $(A) = \text{span}(u_1, \ldots, u_n)$ and the dimension of the range is n.

Exercise 8.17 Let $A = U\Sigma V^T$. Since $U = (u_1, \ldots, u_m)$ is a basis for \mathbb{R}^m we can write

$$b = \sum_{i=1}^{m} (u_i^T b) u_i,$$

Similarly, $V = (v_1, \ldots, v_n)$ is a basis for \mathbb{R}^n so

$$Ax = A(\sum_{i=1}^{n} (v_i^T x)v_i) = \sum_{i=1}^{n} \sigma_i(v_i^T x)u_i.$$

We obtain

$$\|Ax - b\|_{2}^{2} = \|\sum_{i=1}^{n} (\sigma_{i}(v_{i}^{T}x) - (u_{i}^{T}b))u_{i} - \sum_{i=n+1}^{m} (u_{i}^{T}b))u_{i}\|_{2}^{2} = \sum_{i=1}^{n} |\sigma_{i}(v_{i}^{T}x) - (u_{i}^{T}b)|^{2} + \sum_{i=n+1}^{m} |u_{i}^{T}b|^{2} + \sum_{i=n+1}^{m} |u_{i}^{T}b$$

The minimum is obtained for $\sigma_i(v_i^T x) - u_i^T b$ for $i = 1, \ldots, n$, so

$$x = \sum_{i=1}^{n} \frac{u_i^T b}{\sigma_i} v_i.$$

For this particular x we get

$$r = b - Ax = \sum_{i=n+1}^{m} (u_i^T b) u_i.$$

Exercise 8.18 Compute Ax to obtain

$$Ax = A(\sum_{i=1}^{m} \frac{u_i^T b}{\sigma_i} v_i) = \sum_{i=1}^{m} \frac{u_i^T b}{\sigma_i} Av_i = \sum_{i=1}^{m} \frac{u_i^T b}{\sigma_i} \sigma_i u_i = \sum_{i=1}^{m} (u_i^T b) u_i = b,$$

where the last equality holds since $U = (u_1, \ldots, u_m)$ provides an orthogonal basis for \mathbb{R}^m which is the space b belongs to.

Since m < n the matrix has a null space null $(A) = \text{span}(v_{m+1}, \ldots, v_n)$. If x_2 belongs to the nullspace then $A(x + x_2) = Ax = b$ so the solution is not unique. Since the above formula for x does not include a component from the null space it can be characterized as

 $\min \|x\|_2 \text{ such that } Ax = b,$

that is the *minimum norm* solution of the linear system Ax = b.

Exercise 8.19 Since rank(A) = k we note that $\{v_{k+1}, \ldots, v_n\}$ is a basis for null(A) and $\{v_1, \ldots, v_k\}$ is a basis for its orthogonal complement $(\text{null}(A))^{\perp}$. Thus for every x we can write

$$x = x_1 + x_2 = \left(\sum_{i=1}^k c_i v_i\right) + \left(\sum_{i=k+1}^n c_i v_i\right).$$

In order to determine x_1 we compute

$$Ax = A(x_1 + x_2) = Ax_1 + 0 = \sum_{i=1}^k c_i \sigma_i u_i = b = \sum_{i=1}^m (u_i^T b) u_i.$$

Where $(u_i^T b) = 0$, for i = k + 1, ..., m, or a solution doesn't exist. Thus

$$x_1 = \sum_{i=1}^{k} \frac{u_i^T b}{\sigma_i} v_i$$
 and $x_2 = \sum_{i=k+1}^{n} c_i v_i$,

where $c_i, i = k + 1, \ldots, n$, are undetermined parameters.

Exercise 8.20 a) If B has full rank then Bx = 0 if and only if x = 0 so the unique, and only feasible, solution is precisely x = 0.

b) If rank(B) = k < n then B has a non-trivial null space and write $V = (V_k, V_{n-k})$ so that the null space is given by V_{n-k} then the feasible solutions are $x = V_{n-k}c$, $c \in \mathbb{R}^{n-k}$. So in fact we have a regular least squares problem

$$\min_{c \in \mathbb{R}^{n-k}} \| (AV_{n-k})c - b \|_2 \text{ and } x = V_{n-k}c.$$

The above qualifies as a formula. Otherwise continuue and write the normal equations for the above least squares problem.

Exercise 8.21 Let $A = U\Sigma V^T$ and $U = (u_1, \ldots, u_m)$. A solution exists if $b \in \operatorname{range}(A) = \operatorname{span}(u_1, \ldots, u_n)$. We can check this by, for instance, verifying that $u_i^T b = 0$, for $i = n + 1, \ldots, m$. If m >> n it is cheaper to instead check if

$$b - \sum_{i=1}^{n} (u_i^T b) u_i = 0.$$

If we split the matrix $U = (U_1, U_2)$ then the same criteria can be written as $U_2^T b = 0$ or $b - U_1 U_1^T b = 0$.

Exercise 8.22 The normal equations can be derived by the identity

$$\min_{x} \|Ax - b\|_{2}^{2} + \lambda^{2} \|x\|_{2}^{2} = \min_{x} \left\| \begin{pmatrix} Ax - b \\ \lambda x \end{pmatrix} \right\|_{2} = \min_{x} \left\| \begin{pmatrix} A \\ \lambda I \end{pmatrix} x - \begin{pmatrix} b \\ 0 \end{pmatrix} \right\|_{2}$$

The last is a regular least squares problem with an extended matrix. The normal equations are

$$\begin{pmatrix} A^T & \lambda I \end{pmatrix} \begin{pmatrix} A \\ \lambda I \end{pmatrix} x = \begin{pmatrix} A^T & I \end{pmatrix} \begin{pmatrix} b \\ 0 \end{pmatrix} \text{ or } (A^T A + \lambda^2 I) x = A^T b$$

Now we can derive the solution formula using the decomposition $A = U\Sigma V^T$. Since $A^T A + \lambda I = V\Sigma^T \Sigma V^T + \lambda^2 V V^T = V(\Sigma^T \Sigma + \lambda^2 I) V^T$ and $A^T b = V\Sigma U^T b$ we obtain the solution

$$x_{\lambda} = V(\Sigma^T \Sigma + \lambda I)^{-1} \Sigma U^T b = \sum_{i=1}^n \frac{\sigma_i}{\sigma_i^2 + \lambda^2} (u_i^T b) v_i.$$

Too see that the normal equations are not ill-conditioned we look at $A^T A$ which has singular values $\sigma_i^2 + \lambda^2 \ge \lambda^2$. So the addition of the regularization parameter removes the small singular values and makes the condition number smaller.