



# MAT 331: Mathematical problem solving with computers

## Spring 2009

[Home](#)  
[General Information](#)  
[Syllabus and Homework](#)  
[The textbook](#)  
[Classes](#)

### Welcome to MAT 331

This course serves as an introduction to computing for the math student. After a general introduction to the use of the computers, we will turn to more mathematical problems. The emphasis of this course is on the problem solving portion of the title: we will take a series of problems and try to find solutions (or approximate solutions), keeping in mind that we have access to computers. We will discuss the problems and development of necessary mathematics, and then we will turn to the computers to explore and work out the solutions.

*Copyright 2008 Stony Brook University*



# MAT 331: Mathematical problem solving with computers

## Spring 2009

Home  
General Information  
Syllabus and Homework  
The textbook  
Classes

**Place and time:** S-235 Math Tower, Tuesdays and Thursdays 11:20-12:40.

**Instructor:** [Robert Findley](#)

- Office: 3-122 in the Math Tower.
- email: rfindley at math.sunysb.edu
- Office Hours: We 11-12 in MLC, 12-1 in 3-122 of the Math Tower

**Textbook** The **text** for this course is a set of notes currently being written by Professors Sutherland and Simanca. These are available on the class web page; they may be revised somewhat as the semester progresses. You might find it useful to obtain a manual about Maple (a manual can be downloaded from [Maplesoft](#)). Most of what you need will be covered in class, but it is often useful to have a reference at hand.

**Grades policy** Two or three projects, in total 50% of the grade; problems, 20% of the grade; class participation 10% of the grade; Maple test, 20% of the grade.

*Both the expository and computational aspects of the project write-ups will be graded.*

Discussing homework and projects is encouraged, although, *each student will be responsible for turning in a write-up of the problem and solution.* This implies that files cannot be shared in any way, only preliminary general discussion is aloud.

All grades will be posted in [Blackboard](#).

**Projects and Problems** There will be a number of problems assigned, as well as two or three projects.

A problem is a homework assignment, that you should be able to complete in at most a couple of hours. Each problem counts for the same amount of points, whether they are easy or difficult, and there will be between two and four weekly problems. Each problem should contain a description of the steps performed. These problems must be written individually. Identical and too similar submissions will be

penalized (see Academic Integrity).

A project is more like a term paper you will be expected to devote a significant amount of time to doing it, as well as taking care with the presentation. The expository part of each project should contain an introduction, a detailed description of the problem or topic, what means were used in solve it, and the solution. These write-ups should be produced by each student individually, and should be detailed enough so that someone who has not taken the class can read and understand them, and will believe the solution is correct. In particular, all non-standard mathematical terms used should be defined. These write-ups are often acceptable for the mathematics writing requirement. All references must be explicitly quoted. Cutting-and-pasting from the Internet is not aloud. These papers will be submitted through SafeAssign, therefore, violation of these rules will be detected, and penalized (see Academic Integrity).

All problems and projects will be posted in this website.

**Computers** We will use the math computer lab in S-235 of the math tower; this lab contains 30 Apple iMac computers running Ubuntu Linux. While we will be the Linux machines in class, much of the work can be done on other systems. We will rely mostly on Maple (a program that can do algebra, calculus, graphics, etc.), although if other tools are better suited to the task, we may make use of them. No previous experience with computers is needed.

Maple is available for most platforms (Windows, Macintosh, Linux, . . .); Stony Brook students can buy a copy of Maple at the Seawolves Market Place for five dollars. Maple is also available at all Stony Brook Sinc Sites. You might be able to download a copy for Windows or Mac from [Softweb](#).

The mathlab computers can be accessed from off campus via the internet. For remote graphical access to MathLab follow [these](#) instructions.

You can also use ssh to download or upload files from and to the MathLab. In this way, you can download your Maple files to the computers you are working on. The application ssh is installed at all Sinc sites. For ssh instructions, read [here](#).

### **Information for students with disabilities**

If you have a physical, psychological, medical, or learning disability that may impact your course work, please contact Disability Support Services at (631) 632-6748 or

<http://studentaffairs.stonybrook.edu/dss/>. They will determine with you what accommodations are necessary and appropriate. All information and documentation is confidential.

Students who require assistance during emergency evacuation are

encouraged to discuss their needs with their professors and Disability Support Services. For procedures and information go to the following website: <http://www.sunysb.edu/ehs/fire/disabilities.shtml>

**Academic Integrity** You may discuss the assignments in this course with classmates, however *each student's submission must be his or her own work*. Any evidence that a submission is not the student's individual work (this includes the following: the submission has been copied, shared, or transmitted in any way between students, has been taken from the Internet, has been written by others in previous semesters) will be regarded as evidence of academic dishonesty. Additionally, any evidence of sharing of information or using unauthorized information during an examination will also be regarded as evidence of academic dishonesty.

Be advised that any evidence of academic dishonesty will be reported to the Academic Judiciary. Those involved will be prosecuted to the fullest extent permitted by the University and College laws.

Representing another persons work as your own is always wrong. Faculty are required to report any suspected instance of academic dishonesty to the Academic Judiciary. For more comprehensive information on academic integrity, including categories of academic dishonesty, please refer to the academic judiciary website at <http://www.stonybrook.edu/uaa/academicjudiciary/>

**Critical Incident Management** Stony Brook University expects students to respect the rights, privileges, and property of other people. Faculty are required to report to the Office of Judicial Affairs any disruptive behavior that interrupts their ability to teach, compromises the safety of the learning environment, and/or inhibits students ability to learn.

**Classroom etiquette** The following is not allowed during class: cell-phone use (this includes no text-messages), producing non-class work on the computer, side conversations when the instructor or other student is addressing the class. Also, you should try to arrive on time. Let me know if you need to leave early.

**Communication** Many class communications will be done by email. You should check your email once a day. I will use the email address stored in Blackboard. Make sure that this address is up to date.



# MAT 331: Mathematical problem solving with computers

## Spring 2009

[Home](#)  
[General Information](#)  
[Syllabus and Homework](#)  
[The textbook](#)  
[Classes](#)

### Schedule

Homework and projects must be submitted through **Blackboard** (always remove output of the Maple file before uploading.) Below is the list of required homework and projects, in reverse order.

**Homework 9; Due Friday, May 8th.**  
**IFS Notes. Right click to save IFS worksheet.**

Week of Apr 27th. Fractals. **Project 3; Due Sunday, May 10th.**

Week of Apr 20th. Fractals. **Homework 8; Due Wednesday, April 29th**  
**<-- NEW DEADLINE**

Example Problem solutions: **Right click and save.**

Week of Apr 13th. Fractals. Study for Maple test. (on Tuesday, April 21)  
**Example Problems. Syllabus.**

**Project 2; Due Sunday, April 5rd. (Extended)**

Week of Mar 30th. RSA Finale, Primes & Fractals.

Week of Mar 23rd. Number Theory, RSA Encryption. **Homework 7; Due Monday, March 30th.**

**Project Make-up/bonus; Due Wednesday, nMar 25th**

Week of Mar 16th. Cryptography. Number Theory. **Homework 6: Do two of the three problems; Due Monday, March 23rd.**

Week of Mar 9th. Cryptography. Homework 5: Problems 1-4; Due Monday, March 16th.

Week of Mar 2nd. Cryptography.  
**Homework 5; Due Monday, March 16th.**

Week of Feb 25th. Curve Fitting - Cryptography.

**Project 1; due Friday March 6th**

**Project 1 Data.**

**Homework 4: Problems 1, 2, and 3; due Monday March 2nd.**

Week of Feb16th: Curve Fitting. **Homework 3: Problems 9, 10, 12; due**

## Friday Feb 27th

Week of Feb 9th: Curve Fitting. **Homework 2: Problems 6, 7, 8; due Monday Feb 16th.**

Week of Feb 2nd: Preliminaries about Maple. Curve Fitting. **Homework 1: Problems 1, 2 and 3; due Monday Feb 9th.**

Week of Jan 26th: Preliminaries about Maple. If you have a computer, get the Maple CD from Seawolvers Market and install it. Also, submit **Homework 0, due Monday Feb 2nd.**

These are the topics to be covered during the term: Preliminaries about Maple, Curve Fitting, Cryptography and a third topic (which will be chosen later in the term.)

The course text in [html](#), in [pdf](#)  
(written by Scott Sutherland and Santiago Simanca)

[Preliminaries about Maple](#) (also [PDF](#))

[Curve Fitting](#) (also [PDF](#))

[A Glider Model](#) (also [PDF](#))

[Cryptography](#) (also [PDF](#))

[A turtle in a fractal garden](#) (also [PDF](#))



# MAT 331: Mathematical problem solving with computers

## Spring 2009

[Home](#)  
[General Information](#)  
[Syllabus and Homework](#)  
[The textbook](#)  
[Classes](#)

Classes are [here](#)

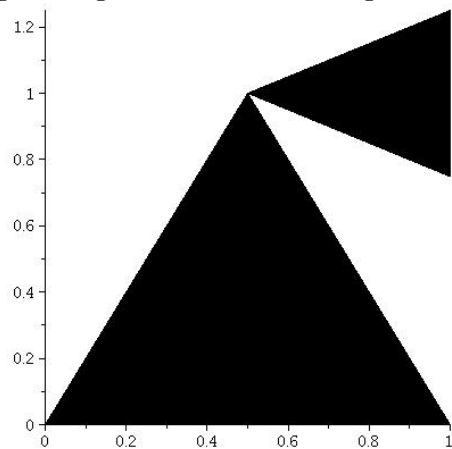
*Copyright 2008 Stony Brook University*



# MAT 331 HW9

Due May 8th, 2009

(1) Find the five parameters for the affine transformation that takes the large triangle to the small triangle:



(2) Create the  $n^{\text{th}}$  approximation of a fractal of your choice using an iterated function system (that we haven't already created in class, of course.)

# MAT 331 Project 3

Due May 10, 2009

Your final project consists of two parts.

(1) Use `TurtleCmd` or an iterated function system to create a fractal of your choice which has box dimension  $d = \log(7)/\log(6)$ .

(2) Write a procedure that takes two parameters,  $n$  and  $a$ , and then draws the Koch approximation  $K_{n,a}$  with recursion depth  $n$  using `TurtleAngle`  $a$ . Compute the box dimension of the Koch curve  $\lim_{n \rightarrow \infty} K_{n,a}$  for at least three different choices of  $a$ . (We've already done the case  $a = 60$  in class!). What do you think happens as  $a$  approaches 90? What about when  $a$  approaches 0?

BONUS: Write an equation  $d(a)$  expressing box dimension as a function of `TurtleAngle`.

As in the previous projects, you will need to attach a text file explaining your Maple file, and containing the relevant calculations and answers to questions from both parts.

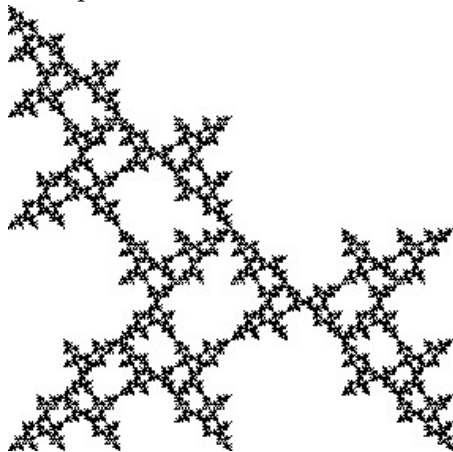
# MAT 331 HW8

Due April 27, 2009

Each of the following problems is worth 10 pts.

(1) Write a recursive function to compute the  $n^{\text{th}}$  Fibonacci number (see section 5.3.1 in the book).

(2) Compute the box dimension of the following figure. Show your work!



(3) Compute the area of the Koch snowflake. You can assume that the largest triangle (formed by the foundation of the three Koch curves) has area 1. You can do this on paper, but explain your reasoning and show your work.

(4) EXTRA CREDIT: Using `TurtleCmd`, draw a random walk. In other words, write a procedure to generate a random string of commands using the letters LRFB, and then display such a string of reasonable length using `TurtleCmd`.

These are examples of the types of questions that will be on the Maple test. You may use the built-in Maple documentation. However, you may not obtain assistance from your classmates or the internet.

1. Plot the function  $y = \sin(x)$  on the interval  $0 < x < 2\pi$ .
2. Solve the following equations for  $x, y, z$ :

$$x + y + z = 7 \quad (1)$$

$$3x + 2y - 4z = 8 \quad (2)$$

$$2x - 17y + 8z = 2 \quad (3)$$

3. Find all of the critical points of the function  $f(x, y) = x^2 + 3xy + y^4$ .
4. Given the matrices  $A, B$ , and the vector  $v$  below, compute the following expressions.

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 5 \\ 7 & 11 & 13 \end{bmatrix} \quad B = \begin{bmatrix} 6 & 2 & 1 \\ 7 & 8 & 0 \\ 4 & 2 & 9 \end{bmatrix} \quad v = \begin{pmatrix} 20 \\ 9 \\ 3 \end{pmatrix}$$

- (a)  $A \cdot B$
  - (b)  $A^{-1}$
  - (c)  $A \cdot v$
5. Create a function  $f(x) = x^2 + 1$ , and use it to expand compute  $f(f(f(f(x))))$ . Create a list containing  $f(n)$  for  $n$  from 1 to 20.
  6. Modular arithmetic/algebra.
    - (a)  $19^{89} \pmod{1024}$
    - (b)  $5x + 10y = 2 \pmod{91}, 3x - 6y = 7 \pmod{91}$ .
  7. Find the 1000<sup>th</sup> prime number.  
Find a 25 digit prime.  
Factor the number "123456789".
  8. Write a procedure or a function that takes a parameter  $\text{tanpt}$  and plots the tangent line to  $y = \sin(x)$  at the point  $(\text{tanpt}, \sin(\text{tanpt}))$ .

9. Write a `for` loop that adds the non-prime numbers between 1 and 1000.
10. Write a procedure that takes `name` as a parameter and responds: “Hi `name`, there are `N` letters in your name”, where `N` is the correct number.
11. Find the location of the first instance of the string of numbers “234” in the digits of  $\pi$ .

For your Maple test, you will be allowed to use the built-in Maple help index. However, you won't be able to work together, or use the internet. You should make sure you are familiar with the following commands.

- Plotting: plot, plot3d, implicitplot
- Calculus: diff, finding the critical points of two variable functions.
- Algebra/Arithmetic: solve, factor, sum, evalf
- Linear Algebra: Matrix, Vector, evalm
- Modular Arithmetic: msolve, mod, modp
- Prime Numbers: nextprime, ithprime, ifactor, phi, chrem
- String Operations: cat, searchtext, print, length, convert, parse
- Lists and sequences: op, map, seq, nops
- Functional Programming: writing procedures, for loops, if..then

# MAT 331 Project 2

Due Friday, April 3rd

In this project, you must research, implement, and write a paper describing a sophisticated encryption method that we have not discussed in class. Two acceptable and recommended topics are the Playfair cipher or ElGamal encryption and signature scheme. Other encryption methods may be researched, but they must be approved by at the latest Friday, March 27th.

Your submitted work should consist of both a maple worksheet and text document. In the worksheet you should include both an enciphering procedure as well as a deciphering procedure, as well as anything else that is relevant to your research topic. The associated document should be treated as a brief but clearly written research paper with two main parts. The first part should discuss the history of your encryption algorithm, mathematically describe the technique, and evaluate its strengths and weaknesses. The second part should explicitly describe your implementation of the method in Maple. You should explain how your procedure works, giving a step by step walkthrough but not describing specific syntax.

The expository portion of this project is more important than in the first project. Specifically, points will be deducted if it is poorly written or improperly formatted. If pressed, I would say it should be about three pages long, though length alone has no bearing upon grading. *You must cite at least three sources for your research, of which at least one must not be a website.*

# MAT 331 HW7

Due March 30th, 2009

- (1) List the encryption method you have chosen for your Project 2, and cite your three sources.
- (2) Twenty-one pirates are dividing their horde of gold doubloons. Since they are a democratic outfit, they first try to divide the coins evenly, but they find there are 19 coins left over. The “discussion” about how to divide the remaining coins results in only 16 pirates still needing to divide the horde (the remaining five went to a place where you can’t bring money or anything else with you). The redivision among 16 pirates leaves 1 coin left over, and three of the pirates make a grab for it. These three find themselves to be missing their hands after this attempt, and the remaining thirteen pirates decide to divide the share among themselves, leaving the handless ones with nothing. Fortunately, the horde divides evenly among the thirteen. What is the minimum number of coins in the horde?
- (3) Use RSA with the modulus  $n = 119$  and the exponent  $e = 7$ , with the 95-character alphabet consisting of the printable ASCII characters to encrypt the word “Yes”. Recall that the alphabet is given by  
`Alphabet:=cat(op(select(IsPrintable, [seq( convert([i],bytes), i=1..255)])))`  
so that  $Y=58$ ,  $e=70$ ,  $s=84$ . Give your encryption as a list of three numbers.



# Project 1 Extra Credit/Make-up

Due Wednesday, March 25rd

As a last concession for any confusion caused by the instructor change, those of you who did the third part of Project 1 incorrectly have an opportunity to make-up it up for 20 of the original 30 points.

In class, we agreed upon two error functions that were acceptable interpretations of the third method in the project description. The second interpretation, which very few people used, was to minimize the sum of the squares of the shortest distances from each point to the line. In other words If a line  $l$  has slope 1, the shortest path from a given point  $p$  to  $l$  is along a line of slope  $-1$ .

To regain points on your project grade, submit through blackboard by Wednesday, March 25th a Maple file fitting your project data to a line using the error function above, as well as a text file describing this method and deriving the error formula used in your worksheet.

To be fair to those who did the third method correctly, I will replace the lowest grade of any of the three methods with the score from your bonus submission, if the resulting score is higher. Additionally, the project grades will not be curved, so if you choose not to submit a make-up your grade will not be affected.

# MAT 331 HW5

Due Monday, March 23, 2009

Do two of the following problems.

**(1)** (a) Modify the procedure 'OneTimePad' to take a text string as a key. Specifically, have your procedure convert a text string of equal length to your message and use it to encrypt the message. This is sometimes called an Augustus Cypher.

(b) Describe a possible method of cracking a lengthy portion of cypher-text encyphered with the above method. Be specific.

**(2)** (a) If we compose a Vignère Cypher of length 3 and a Vignère cypher of length 5, what type of cypher do we have? Although we begin with eight unknowns, how long is the resulting key? What if our keys were of length 2 and 6?

(b) Write a procedure to that takes two keys of arbitrary lengths  $n$  and  $m$  and applies the Vignère cypher with each of these keys.

**(3)** (a) Write a procedure `AffineMatDecode` that takes as its input an encrypted message, a  $2 \times 2$  matrix  $A$ , and a 2-vector  $B$ , and *decyphers* the message under the assumption that it was encoded with the affine cypher  $x \mapsto A \cdot x + B$ . You may use as a start our existing `AffineMatEncode` procedure. Be careful to have your procedure work independent of the specific alphabet *AlphaNumeric*. You may assume that the affine cypher is decryptable. That is, that the determinant of  $A$  is invertible modulo the length of your alphabet.

(b) Use this to decrypt the following message, which was encrypted with the affine cypher  $A = \begin{bmatrix} 17 & 22 \\ 4 & 53 \end{bmatrix}$ ,  $B = [19, 82]$  using the 95 character alphabet we have often used in class.

15Ti\"tEbTiN7H05.Fcvq?K^e~ [%BFckW

# Homework 5

Due Monday, March 16, 2009

For the following problems, use the 27 letter alphabet:

‘ ‘ abcdefghijklmnopqrstuvwxyz ’ ’

(there is a space in the 0<sup>th</sup> position).

(1) Using a *brute force* method (trial and error), decrypt the message

lbhnpjzznasisentwarn s

which was encoded with a Caesar cypher and the 27 letter alphabet.

*Hint:* use a loop.

(2) Suppose you intercept the message

tfqlvddmvleokovt xkthprd

and you know that it is encoded with a Vignère cypher of length 3, using the above alphabet. Suppose further that you know the message starts with “it ”. Decode it.

(3) The string below was encrypted using an affine cipher on the same 27 letter alphabet. Decrypt it.

fmw segjaweouanerj a ceyqrype aswaheoaqbrqabeafrua eeaojerf afmjeayperjpu

*Hint:* this phrase follows the the typical pattern in English where there are (almost) as many spaces as words (and so spaces are very common), and the letter “e” is also very common. You can use the technique described in chapter 4 of the notes, section 7.3.

(4) (a) Compare and contrast the Vignère cypher of length two to the affine cypher. Both have two unknowns, but which is easier to decrypt?

(b) Now compare the Vignère of length six with the  $2 \times 2$  affine matrix cypher. Both have 6 unknowns. Does the situation change?

*You can type your answer to this problem on your Maple worksheet, or hand it in before class on Tuesday, March 17.*

## MAT 331, Spring 2009

### Project 1: Fitting a line to data

In this project, you are to fit a line to some data in three different ways. First, by using the “usual” least squares method, which minimizes the sums of the squares of the vertical distances between the data and the line. Then, find the line given by assuming that the  $x$  variable may contain the errors, and minimizing the sums of squares of the *horizontal* distances. Finally, you are to find the line that minimizes the sum of the *absolute* distances from the points to the line. You should, of course, also have Maple plot the three lines you find along with the data points.

Your project should consist in two parts: a Maple (or other program) with computations and a paper which explains clearly each step you take and why. In particular, you *must* give a careful derivation of the functions you minimize—you cannot merely quote the formula for the distance from a point to a line. You can use built-in Maple commands to check your results, but you have to perform step by step calculations. You may, of course, assume that the formula for the distance between two points in the plane is known to the reader. Your paper should be written so that someone who has a good working knowledge of mathematics but has never taken this class will understand what you are doing and why.

You should use Maple (or any another appropriate program of your preference) to do the actual computations. *Do not* include false starts, mistakes, or irrelevant calculations in your finished product. Pay attention to spelling, grammar, and sentence structure. **The expository part of this project counts as much as the actual calculations**, and should not be ignored.

Both, the paper and Maple file should be submitted through Blackboard. To submit the paper you must use SafeAssign and to submit Maple file, you must use a “regular” Assignment (as in the homework) You will find a file containing the data points at

<http://www.math.sunysb.edu/~moira/mat331-spr09/data/project1data/yourLastname.txt>,

where *yourLastname* is the your lastname (!).

## Homework 4

Due March 2, 2009

**(1)** (# 11 from previous sheet) Following the book (Fitting a Circle) of the notes, prove that if we describe the circle of center  $(a, b)$  and radius  $r$  using the parameters  $(a, b, k)$ , with  $k = a^2 + b^2 - r^2$ , rather than the more natural parameters  $(a, b, r)$ , then the error function  $H(a, b, k) = E(a, b, \sqrt{a^2 + b^2 - k})$  is quadratic in  $a, b$  and  $k$ . What does this imply about the number of critical points?

**(2)** The following is real data illustrating the number of stores carrying a product (x-axis) vs sales of that product (y-axis). Fit the the data to three curves: one linear, one quadratic, and one exponential. You may use any curve fitting procedures used in class. In one or two sentences, explain which curve you think fits the data better, and suggest a reason why this curve might model the real-world scenario illustrated by the data.

data=[[4056.29, 9667.91], [3572.48, 5249.2], [3720.03, 6792.28], [1826.62, 1604.18], [2770.98, 4096.54], [2155.19, 2251.57], [2250.6, 2177.33], [3394.23, 5654.73], [3783.34, 8140.33], [3577.43, 4904.89], [2044.14, 1550.97], [3054.78, 3174.18], [844.56, 1568.49], [2175.06, 1855.8], [2806.39, 3202.34], [1639.95, 1402.94], [2455.29, 1735.76], [1697.39, 1664.51], [1071.2, 984], [1484.8, 1376.29], [1157.56, 1264.41], [839.34, 1250.43], [641.07, 753.53], [3686.87, 6353.6], [969.75, 1294.29]];

**(3)** A cryptography question will be added here.

## MAT 331, Spring 2009, Problems

1. Use **Maple** to write  $x^5 - 6x^3 - 5x - 3x^4 + 18x^2 + 15$  as a product of *exact* linear factors. By exact, I mean you should leave any non-rational factors expressed as radicals; do not approximate terms like  $\sqrt{3}$  as 1.73205, etc.
2. Consider the planar curve  $\gamma$  defined by  $x^2y^3 + y^2 + y - 2e^x = 0$ . Using **only Maple**, find the slope of the tangent line to the curve at  $(0, 1)$ . Then plot the curve and the tangent line on the same graph.  
*Hint: you might want to use `implicitplot` from the library `plots`. You might find `implicitdiff` helpful, too.*
3. Plot the function  $f(x) = 2\sin x - x^3 - 5$ , for  $x \in [-4, 4]$ . Find all the zeros of the function with an accuracy of 20 decimal digits. *Hint: See `Digits`, `fsolve`, and/or `evalf`. Make sure you find *all the zeroes*. Justify your answer.*
4. Define a **Maple** function  $g$  that, given a positive integer  $k$  yields the sum of the first  $k$  primes. What is  $k$  such that  $g(k) \leq 100,000$  but  $g(k+1) > 100,000$ ? You might find `sum` and `ithprime` helpful.
5. The Fibonacci sequence is a sequence of positive integers defined by recurrence as follows:  $F_0 = 1, F_1 = 1$  and for each integer  $i$  larger than one,  $F_i = F_{i-1} + F_{i-2}$ . Use Maple to find the first 100 terms of the sequence without computing these terms one by one (Hint: See Fibonacci). Use the sequence you found to compute a new sequence  $F_i/F_{i-1}$  and try to guess what is the limit of that new sequence. EXTRA CREDIT: Prove the sequence  $F_i/F_{i-1}$  converges to the limit you guessed (You may need to refine your guess to do so).
6. Find a polynomial of degree 3 passing through the points  $[-3.3, 1], [-2, 2], [0, 3], [1.1, 4]$  by two methods: First solve the equations (as in the book). Second, use the build-in **Maple** command. Check that the two solutions agree.
7. (a) Find two polynomials of degree 5 passing through the points  
 $[-3.3, 1], [-2, 2], [0, 3], [1.1, 4]$ .  
  
(b) Plot the two polynomials and the points in the same graph.  
(c) Is it possible to find more polynomials of degree 5 passing through the 4 points? Why?  
(d) How many polynomials of degree 3 passing through the 4 points can you find?
8. Consider the set of points

$$\begin{aligned} & [ -3, -9/4 ], [ -2, -3/2 ], [ -1, -3/4 ], [ 0, 0 ], [ 1, 3/4 ], [ 2, 3/2 ], \\ & [ 3, 9/4 ], [ 4, 3 ], [ 5, 15/4 ], [ 6, 9/2 ], [ 7, 5.35 ], [ 8, 6 ]. \end{aligned}$$

- (a) Find a polynomial passing through all the points.
  - (b) find a piecewise linear function passing through all the points.
  - (c) Fit the points to a line using the least square method.
  - (d) Plot the three results (the to of the previous problem, and the one you just obtained and the points in the same graph.
  - (e) Compare the results. Which one do you think is best and why?
9. Fit the points  $(-1.9, -4.7), (-0.8, 1.2), (0.1, 2.8), (1.4, -1.2), (1.8, -3.5)$  by means of a quadratic function  $f(x) = ax^2 + bx + c$ , using the least square method. First, do this step by step, as we did in class; then, use the built-in **Maple** command, described in the notes. Check that the two solutions agree.
10. Fit the set of points

$$(1.02, -4.30), (1.00, -2.12), (0.99, 0.52), (1.03, 2.51), (1.00, 3.34), (1.02, 5.30)$$

- with a line, using the least square method we used in class. You will see that this is not a good fit. Think of a better way to do the fit and use **Maple** to do it. Explain in your solution why you think your better way is better.
11. Following the book, prove that if we describe the circle of center  $(a, b)$  and radius  $r$  using the parameters  $(a, b, k)$ , with  $k = a^2 + b^2 - r^2$ , rather than the more natural parameters  $(a, b, r)$ , then the error function  $H(a, b, k) = E(a, b, \sqrt{a^2 + b^2 - k})$  is quadratic in  $a, b$  and  $k$ . What does this imply about the number of critical points?(Note: You do not need **Maple** in this exercise, although you can use it as a word processor. If you're more comfortable with paper, you can turn in a paper instead of a **Maple** worksheet)
12. Fit the points  $[[1, 1], [2, 3], [3, 1], [1, -1.1]]$  to a circle.
13. Following the on-line notes, find a set of 25 points which are "approximately" in a cubic and:
- (a) Find the cubic that best fit the points by the least squares method using the CurveFitting package.
  - (b) Find the cubic that best fit the points by the least squares method without using the CurveFitting package.
  - (c) (EXTRA CREDIT) Find the corresponding spline.
14. *[In this problem use **Maple** only as a word processor. If you're more comfortable with paper, you can turn in a paper instead of a **Maple** worksheet.]* Let  $n$  points of the form  $(r_i, r_i^2)$ ,  $i = 1, 2, \dots, n$ , be given. What is the quadratic function  $f(x) = ax^2 + bx + c$  that best fits them? **Prove** your answer. Does it depend on the optimization method (least square or others)?

Send me (moira@math.sunysb.edu) with your name, mathematical and computer background, hopes and fears about this course and an explanation of what are you planning to do with your major. Please include any other information about yourself that you think it could be useful for me. Finally, include an acknowledgment of the academic dishonesty policy in this course (you can find it in the syllabus, in the webpage of this course). Please write in the subject of your email message (and of all your email messages to me) MAT331. This will contribute to organize our email communication. Make sure you send this message. Future homeworks will be graded only if I received it.



# Preliminaries about Maple

Maple is a comprehensive general purpose computer algebra system which can do symbolic and numerical calculations and has facilities for 2- and 3-dimensional graphical output. Calculus courses may be structured so that Maple can be used as a tool that will help you gain a complete understanding of the material discussed. You will discover that the capability of Maple goes well beyond the realm of calculus; it is a tool that is used more and more in education and scientific research in mathematics and engineering.

This chapter is a brief introduction to the use of Maple in the environment that you will discover at Stony Brook. It is not exhaustive and its sole purpose is to guide you through the first steps needed to start using this tool. Once you become familiar with the basics of Maple and can appreciate its advantages, we encourage you to experiment with it, not only to learn mathematics in more detail but also to help you in other courses and your own research. More advanced documentation is available on line, in the computer lab, and in the library.

Maple runs on many different types of computers and operating systems. For example, you might use Maple from a Sun workstation running Solaris (a flavor of Unix) in the computer lab, from a Windows computer in the library, and from a Macintosh at home. Maple behaves similarly on all of them; we will usually point it out when we use something specific to one type of computer.

---

## Subsections

- [Starting a Maple Session](#)
- [Basic Maple](#)
  - [Maple for numeric calculations](#)
  - [Maple for symbolic manipulations](#)
  - [Graphing in Maple](#)
- [The Maple worksheet](#)
  - [Introducing the worksheet](#)
  - [Worksheet basics: the front-end and the kernel](#)
  - [Online Maple help](#)
  - [Documenting and structuring your worksheet](#)
- [Assignments, Functions and Constants](#)
  - [Assignment statements](#)
  - [Variables and subsequent assignments: unassignments](#)
  - [Functions known to Maple. How to define your own functions](#)
  - [Special constants and reserved names](#)
- [The limit command](#)
- [The diff \(and Diff\) command](#)
  - [Higher order derivatives](#)
  - [Implicit differentiation](#)

- [The int \(and Int\) command](#)
  - [Symbolic integrals](#)
  - [Numerical integration](#)
  - [Approximations through Riemann Sums](#)
  - [Multiple integrals](#)
- [The subs command](#)
  - [Why is subs useful?](#)
- [Plotting with Maple](#)
  - [Plotting several functions \(or curves\) on the same axes](#)
  - [Fancier plotting](#)
- [Exercises](#)

---

[next](#) [up](#) [previous](#)

**Next:** [Starting a Maple Session](#) **Up:** [Using Maple](#) **Previous:** [Using Maple](#)

*Translated from LaTeX by Scott Sutherland  
2002-08-29*



# Chapter 1

## Preliminaries about Maple

Maple is a comprehensive general purpose computer algebra system which can do symbolic and numerical calculations and has facilities for 2- and 3-dimensional graphical output. Calculus courses may be structured so that Maple can be used as a tool that will help you gain a complete understanding of the material discussed. You will discover that the capability of Maple goes well beyond the realm of calculus; it is a tool that is used more and more in education and scientific research in mathematics and engineering.

This chapter is a brief introduction to the use of Maple in the environment that you will discover at Stony Brook. It is not exhaustive and its sole purpose is to guide you through the first steps needed to start using this tool. Once you become familiar with the basics of Maple and can appreciate its advantages, we encourage you to experiment with it, not only to learn mathematics in more detail but also to help you in other courses and your own research. More advanced documentation is available on line, in the computer lab, and in the library.

Maple runs on many different types of computers and operating systems. For example, you might use Maple from a Sun workstation running Solaris (a flavor of Unix) in the computer lab, from a Windows computer in the library, and from a Macintosh at home. Maple behaves similarly on all of them; we will usually point it out when we use something specific to one type of computer.

### 1 Starting a Maple Session

How to start a Maple session varies somewhat in different environments. On the Windows machines, you can usually either start Maple by selecting the option from the start menu, or by finding the appropriate icon (it usually looks like a Maple leaf) and clicking the mouse on it. In a Unix environment, there may be an icon or menu item for Maple; you can usually always just issue the command `xmaple &` or `maple -x &` to start it.

After a brief pause, you should see a window that looks similar to that in Figure 1; this is Maple's graphical (or worksheet) interface, where you will type your commands and see the

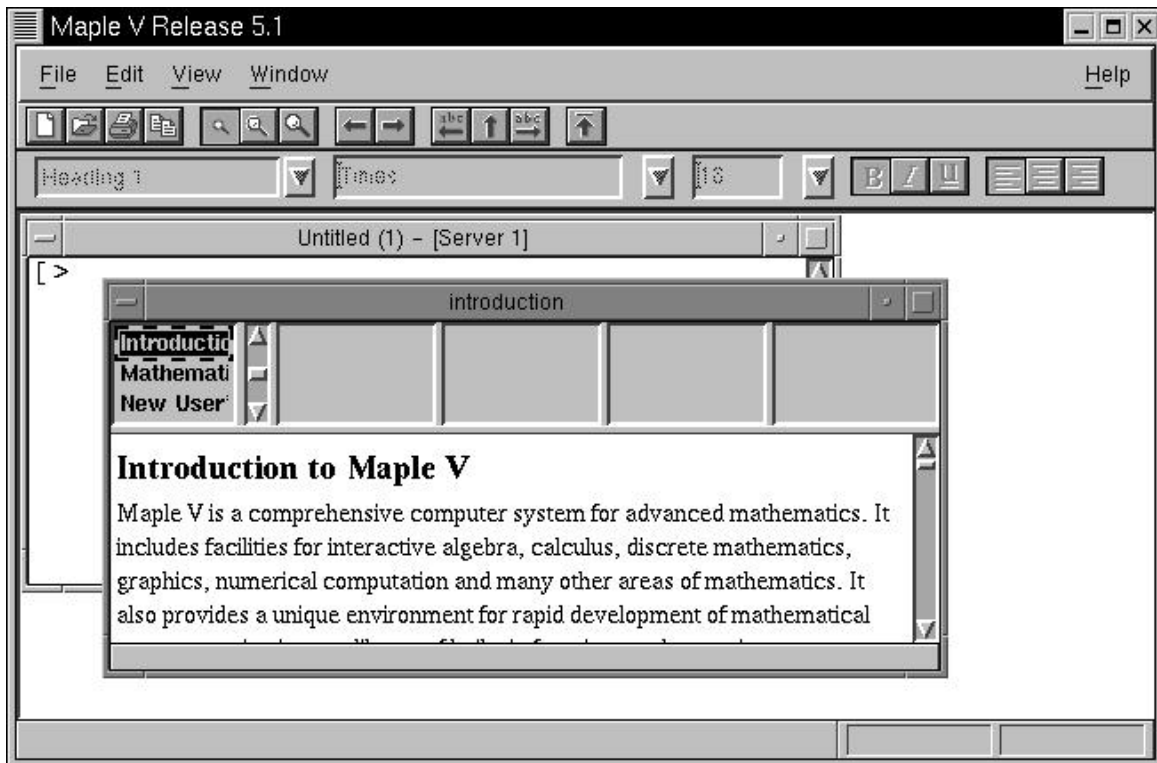
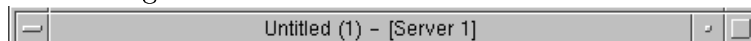



Figure 1: A newly-opened Maple worksheet.


results. We will go into more details of the worksheet interface in §3. As in the figure, your session may have started with two windows open inside the main Maple window. The one in front in Figure 1 is a help window; we'll come back to this in §3.3. But for now, let's concentrate on a worksheet where we can enter some Maple commands. The top part of this window should look something<sup>1</sup> like this:



This window is a Maple worksheet, and is a place to enter your Maple commands. For now, click on the button which maximizes the current worksheet; this is on right side of the worksheet's title bar, and looks like .<sup>2</sup> Clicking the mouse on this should make the worksheet fill most of the Maple window.

Most users typically use the Maple's graphical, or worksheet, interface; this is what starts you add the `-x` to the `maple` command in Unix. However, there is also a text-based interface, which you might see if you just type `maple`, omitting the `-x`.

<sup>1</sup>This particular image was made on a computer running Unix. Under Microsoft Windows or on a Macintosh, the symbols around at the edges will be slightly different, but with similar functions.

<sup>2</sup>Again, if you are running another type of computer, the button may look a little different. For example, under MS-Windows, it looks similar to this: .

```

|\~/|      Maple V Release 5.1 (SUNY Stony Brook)
._|\|    |/|_ . Copyright (c) 1981-1998 by Waterloo Maple Inc. All rights
\  MAPLE / reserved. Maple and Maple V are registered trademarks of
<_--- _---> Waterloo Maple Inc.
      |      Type ? for help.
>

```

This interface is useful when you want to run Maple in a non-graphic environment; for example, over a dial-up connection. However, keep in mind that while all the usual Maple commands work, this interface is much more limited—it is difficult to edit your commands or save your work, and graphic commands such as `plot` are much more limited. If you get into this mode by mistake, you can just type the Maple command `quit`; to exit.

## 2 Basic Maple

You can use Maple for numerical calculations, for symbolic manipulations, and for graphing purposes. *All Maple commands must be ended with a semicolon ; and the corresponding stroke of the return key.* To familiarize yourself with it, try to run the following examples and analyze the results. The commands that you type in and the Maple results are shown here in a way that is similar to what you see on the screen.

### 2.1 Maple for numeric calculations

Most of the time you will be using Maple as a calculator. You can also write programs in Maple, but usually you will be using the tool expecting an immediate response. You type in your statement at the Maple prompt `>` and Maple responds with an answer:

```
> 1/2+3;
```

$$\frac{7}{2}$$

The symbols for the basic operations are `+`, `-`, `*` and `/`. Exponentiation is denoted by `^`:

```
> 77*(3^5+5^2)/(99-47);
```

$$\frac{5159}{13}$$

In the first example, the inverse of 2 is added to 3. In the second, we input  $\frac{77(3^5 + 5^2)}{99 - 47}$ , and Maple computes its value before printing the output.

You must always use the asterisk `*` for multiplication. Forgetting it in expressions like `sin(2*x)` or `4*x` is easy, and may result in a syntax error message. Even if you don't get a syntax error, Maple will not treat the expression as you expect.

Maple works primarily with whole numbers or fractions. However, you can force it to write the decimal expansion of any number with `evalf`:

```
> evalf(100*Pi);
```

314.1592654

```
> b:=sqrt(2);
```

$b := \sqrt{2}$

```
> evalf(b);
```

1.414213562

```
> evalf(b,47);
```

1.4142135623730950488016887242096980785696718754

Be careful with this example. The *assignment* of the variable `b` to  $\sqrt{2}$  is done with `:=` rather than with `=`. This will be the case of any other assignment operation you do in Maple, as we discuss below. The equals sign `=` is used to signify that two things are equal, not to set one to the other.

Notice the difference in the following two examples:

```
> sum(1/(2*i),i=1..10);
```

$\frac{7381}{5040}$

```
> sum(1.0/(2*i),i=1..10);
```

1.464484127

In both cases, we compute  $\sum_{i=1}^{10} (2i)^{-1}$ . However, in the second case, the presence of `1.0` indicates to Maple the terms in this computation are approximate, and gives you a decimal expansion of the result.

## 2.2 Maple for symbolic manipulations

We can ask Maple to give us a formula for the sum of the cube of the first  $n$  integers:

```
> sum(i^3,i=1..n);
```

$$\frac{1}{4}(n+1)^4 - \frac{1}{2}(n+1)^3 + \frac{1}{4}(n+1)^2$$

We can ask to factor the resulting polynomial in  $n$ :

```
> factor(%);
```

$$\frac{1}{4}n^2(n+1)^2$$

In this example, Maple's "ditto operator" % is used<sup>3</sup> to refer to the result previously obtained. Thus, in effect, the example above is equivalent to

```
> factor(1/4*(n+1)^4 - 1/2*(n+1)^3 + 1/4*(n+1)^2);
```

$$\frac{1}{4}n^2(n+1)^2$$

One can expand or factor conveniently:

```
> p:=(x+1)^4*(x-6)^3;
```

$$p := (x+1)^4(x-6)^3$$

```
> expand(%);
```

$$x^7 - 14x^6 + 42x^5 + 112x^4 - 287x^3 - 882x^2 - 756x - 216$$

```
> factor(%);
```

$$(x+1)^4(x-6)^3$$

We can use Maple to solve equations (if  $b$  still has the value  $\sqrt{2}$  from the previous section, first *unassign* it by executing the command `b:='b'`):

---

<sup>3</sup>In Maple V release 4 and earlier, Maple used the double-quote " to refer to the result of the previous command instead of %. If you encounter an older Maple program, you might have to change this, among some other things.



```
> quadeq:=a*x^2+b*x+c;
```

$$quadeq := ax^2 + bx + c$$

```
> solve(quadeq=0,x);
```

$$\frac{1 - \sqrt{2} + \sqrt{2 - 4ac}}{2a}, \frac{1 - \sqrt{2} - \sqrt{2 - 4ac}}{2a}$$

Notice that `solve` provides both roots of the quadratic polynomial.

Maple also divides polynomials. Here the command you execute has a syntax which is a little bit more complicated. Consider the following example:<sup>4</sup>

```
> p:=x^10-1;
   q:=x-1;
   divide(p,q,b);
```

```
p := x10 - 1
q := x - 1
true
```

We assign the polynomials  $x^{10} - 1$  and  $x - 1$  to the variables `p` and `q`, respectively, and then ask Maple to divide `p` by `q`.<sup>5</sup> Maple replies with the statement `true` to indicate that the polynomial  $q = x - 1$  divides the polynomial  $p = x^{10} - 1$  exactly, and places the quotient in the third argument of the `divide` command. Thus, if you want to know the value of the quotient, you must see what the variable `b` stands for now:

```
> b;
```

$$x^9 + x^8 + x^7 + x^6 + x^5 + x^4 + x^3 + x^2 + x + 1$$

Of course, the answer produced by `divide` will be `false` if the polynomials do not divide each other exactly. Then one can use the commands `quo` and `rem` to find the *quotient* and *remainder* of the division:

```
> r:=x^2+2*x+3;
```

---

<sup>4</sup>In the worksheet interface, we can insert linebreaks by holding the *shift* key and the *return* key simultaneously. Maple does not execute what we type until we hit the return key unshifted.

<sup>5</sup>Of course, we could have done all of this in a single Maple command `divide(x^10-1, x-1, b)`; if we had no intention of using `p` or `q` again.

$$r := x^2 + 2x + 3$$

```
> divide(p,r,c);
```

*false*

```
> quo(p,r,x);
```

$$x^8 - 2x^7 + x^6 + 4x^5 - 11x^4 + 10x^3 + 13x^2 - 56x + 73$$

```
> rem(p,r,x);
```

$$-220 + 22x$$

We would have obtained an error if we would have tried `divide(p,r,b)`, trying to place the quotient of the division in the variable `b`: it is no longer a variable since it has a value from the previous example:

```
> divide(p,r,b);
```

Error, wrong number (or type) of parameters in function divide

The above command is equivalent to `divide(p,r,x^9+x^8+x^7+x^6+x^5+x^4+x^3+x^2+x+1)`, which makes no sense at all. If we insist on calling the result of the division `b`, it would have been necessary to *unassign* the value of `b` first:

```
> b:='b';
```

$$b := b$$

We could also have done this within the command itself, using `divide(p,r,'b')`.<sup>6</sup> Further discussion on this is given in section 4.2.

We can use `solve` for systems of equations. For example:

```
> solve({1*x+(1/2)*y+(1/6)*z =1,
        (1/2)*x+(1/6)*y+(1/24)*z =0,
        (1/6)*x+(1/24)*y+(1/120)*z =1}, {x,y,z});
```

---

<sup>6</sup>In both cases, it is very important to use the proper quotation marks; using `"` or `'` means a very different thing, and would not work. The quotation mark used here is the 'single quote' or 'right quote', and on most keyboards is on the same key as the double quote, near the enter key.

$$\{x = 39, y = -216, z = 420\}$$

Maple also solves inequalities:

```
> solve(x^2-4*x-7>0,x);
```

$$\text{RealRange}(-\infty, \text{Open}(2 - \sqrt{11})), \text{RealRange}(\text{Open}(2 + \sqrt{11}), \infty)$$

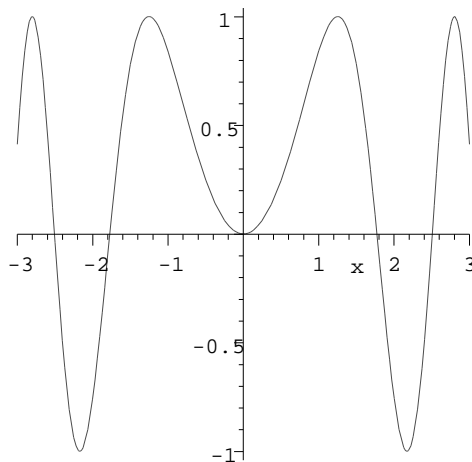
```
> solve(abs(1-x^2)<1/2,x);
```

$$\text{RealRange}\left(\text{Open}\left(-\frac{\sqrt{6}}{2}\right), \text{Open}\left(-\frac{\sqrt{2}}{2}\right)\right), \text{RealRange}\left(\text{Open}\left(\frac{\sqrt{2}}{2}\right), \text{Open}\left(\frac{\sqrt{6}}{2}\right)\right)$$

## 2.3 Graphing in Maple

Let us now use the graphics facilities of Maple. It is quite straightforward to draw the graph of a simple expression. For example, to graph the function  $\sin(x^2)$  over the range  $-3 \leq x \leq 3$ , we do the following.

```
> plot(sin(x^2), x=-3..3);
```



Notice that within the worksheet interface, clicking the right mouse button on the displayed graph causes a menu to pop up that allows you to change some of its properties. Clicking the left mouse button causes a box to appear around the graph and the buttons on the toolbar to change. You can use the box to resize the graph, and the buttons to change attributes such as the style of the axes and the aspect ratio.

In section 9 we shall discuss in further detail the graphing features of Maple.

## 3 The Maple worksheet

Maple's graphical interface allows you to keep a collection of commands and their results in one place, to save your work, and resume previous sessions. You can also add comments and text to your worksheet, to make it more readable and to describe the process of solving the problem. It also enables you to structure your worksheet into sections of related items.

### 3.1 Introducing the worksheet

As we have seen, you enter Maple commands at the `>` prompt, and the results appear in your worksheet. If you go back and change a Maple command and re-execute it, the result in the worksheet changes.

At the top of the screen is a menu bar (with entries like **File**, **Edit**, and so on): clicking on each of these words gives you a menu which allows you to affect your session in various ways. For example, clicking on the **File** menu and selecting **Save** will allow you to save your current worksheet for later use; selecting **Open** from the same menu allows you to load an existing worksheet.

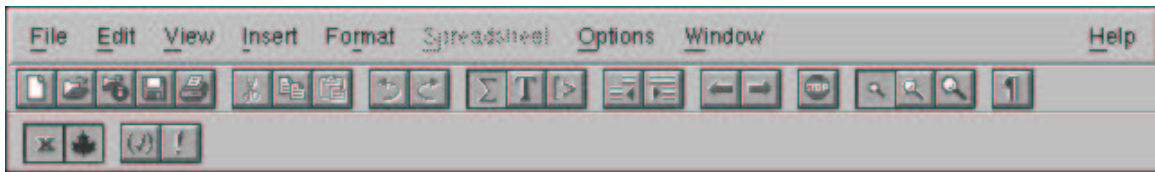


Figure 2: The menu bar, button bar, and context bar.

Below the menu bar is a collection of buttons which are short-cuts for items found on the various menus. Clicking on one of these is much quicker than searching for the command on the proper menu. For example, the button that looks like a floppy disk (📁) saves the current worksheet, and the magnifying glasses (🔍) adjust the zoom factor (or magnification) of the current worksheet. The stop sign interrupts a computation, which is very useful if you ask Maple to do something that takes a very long time, such as to find the factors of  $2^{1000} - 1$ .

Below this is the “context bar”, which changes depending on what kind of item you currently have selected. For example, as mentioned in §2.3, when graphics are selected there are buttons to adjust the style of graphics (points or lines for 2-dimensional graphics, shading on 3-dimensional graphics), axes style, and so on. When a text item or comment section is chosen, selections for typeface, centering, and so on appear.

### 3.2 Worksheet basics: the front-end and the kernel

It is important to remember that the worksheet is only a record of a Maple session, not the session itself. Since you have the ability to edit and reorder it, it can sometimes appear to

contain inconsistent information. Here is a small example:

First, suppose we issue the following commands:

```
> one := sin(x)^2 + cos(x)^2;
```

$$one := \sin(x)^2 + \cos(x)^2$$

```
> diff(one,x);
```

0

The command `diff(one,x)` asks Maple to compute the derivative of the expression called `one` with respect to the variable `x`. The result is, of course, 0 because  $\sin^2 x + \cos^2 x$  is a constant.

Now, we go back and change the `+` in the definition of `one` to a `-`, and execute the statement. However, we do not execute the `diff` command. Our worksheet now looks like this:

```
> one := sin(x)^2 - cos(x)^2;
```

$$one := \sin(x)^2 - \cos(x)^2$$

```
> diff(one,x);
```

0

This seems to contain an error (the derivative should be  $4 \sin(x) \cos(x)$ ). Obviously, this is because we we didn't execute all the statements. If you place your cursor on the line with `diff(one,x)` and hit the return key, the worksheet will be correct again.

This may seem like a silly discussion, but this apparently obvious example is to make an important point: Maple only “sees” your commands in the order you execute them, *not* in the order you see them in the worksheet. Also, if you load a worksheet from a previous session (or made by someone else), any assignments and results will not be accessible to the Maple session until you execute those statements, *even though the results appear on your screen*.

To help you understand why this is, we point out that a Maple session typically consists of two processes (or programs) which communicate with each other: the Maple kernel and the “front-end”, or user interface. The front-end is the part that interacts with the user, accepting your commands, showing the results of your computations, and so on. The kernel is the part that does the actual computations. In certain situations (for example, when using a “parallel Maple kernel”), it is possible to have more than one Maple kernel associated with the same front-end.

### 3.3 Online Maple help

Maple has an extensive and very useful online help system, containing information on all of the commands, as well as a number of tutorials and example worksheets.<sup>7</sup> Spending some time exploring the help system is well worth the effort. There are a few ways to get into the system; one obvious one is to select from the **Help** menu (on the right of the menu bar). But if you know the specific item you want help on, you can also do this in a quicker way.

Suppose you want information on Maple’s **factor** command. Enter the Maple command  
`> ?factor`

(here the semicolon is optional). Then Maple creates another window that very likely will contain more information than you really need or care about. You can read all or part of the information displayed in this window by moving the text down with the *scroll bar* to the right of the window. You can also click in the help window, and use the arrow keys on your keyboard.

You should notice that at the end of the help window you will find a good set of examples on the uses of the command you are inquiring about. This examples can be copied and pasted in your worksheet and executed. In this way, you can see what the command really does.

Usually there is also a section labeled “See Also”, which points you to related commands. Each of these words is a link to another help page. It works like a web browser: clicking the underlined word takes you to the help page on that topic.<sup>8</sup> You can use the back button (☐) to return to the previous help page.

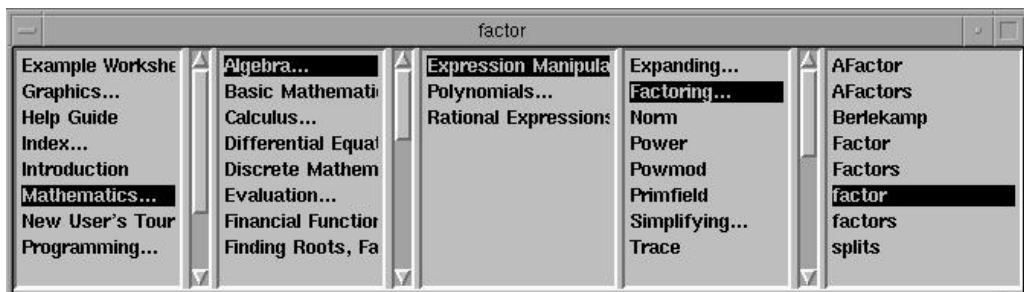


Figure 3: Maple’s help index.

At the top of each help page is the help index (see Fig 3), which is a table of contents for the entire help system, organized by topic. You can use this to browse for commands you might not otherwise have known about.


Finally, on the **Help** menu are entries called “topic search” and “full text search”. The first looks for help on a specific topic, and the latter searches the text of every page in the help system for the word you ask. Both of these can be valuable resources when used properly.


<sup>7</sup>There is even help on the help system. Select **Using Help** from the **Help** menu, or give the command `?helpguide`.

<sup>8</sup>If you prefer, holding the **Control** key while clicking on the link opens the page in a new help window, allowing you to view both at the same time.

### 3.4 Documenting and structuring your worksheet

Aside from the great advantage of making it easy to keep track of your commands and their results, the worksheet has some features of a word processor, allowing you to insert comments and text which describes your solution (or anything else you like). This can be very important, especially if you intend to show your worksheet to someone else (or even look at it yourself later).

To enter text into a worksheet, click on where you'd like it to go, and hit the text-mode button (it looks like )<sup>9</sup>. Note that this button will convert a Maple input region (that is, one with a > prompt) to a text region—you'd probably prefer to *insert* your comments rather than change your commands to comments. To do this, either click before the > prompt, or select **Paragraph...Before** from the **Insert** menu.

In the text region, you can type whatever you like, using the context bar to adjust formatting such as the typeface, centering, etc. You can even include mathematical expressions in your text: click the “math input” button ()<sup>9</sup> and type the Maple expression for the math you want. When you are done, use the text-mode button to return to entering text. You can also paste the output from a previously executed Maple command into the text area.

The worksheet interface also allows you to group parts of your worksheet into sections. These sections can contain paragraphs and Maple execution groups, and subsections. You can recognize a section by a large square bracket to the left, with a box at the top of it. Clicking on the box collapses the section, showing only the its title. Collapsing all sections allows you to see the structure of your document in outline form. To expand a collapsed section, click on the box next to the title.

For an extensive description of how to document and structure your worksheets, see the “Documenting your Work” section in the online help. This can be invoked with `?documenting`.

## 4 Assignments, Functions and Constants

Let us discuss the syntax of some basic Maple commands, as well as some functions and constants that Maple knows.

### 4.1 Assignment statements

The assignment is the most basic and used of Maple statements. It allows you to give names to expressions you want the computer to remember for later use. It has the *colon-equals* format `name := expression`. The thing on the left side of `:=` is a *name*, i.e. a sequence of letters and digits that begins with a letter. Certain names must be avoided because Maple uses them for special purposes (see section 4.4):

---

<sup>9</sup>Or select **Math Input** from the **Insert** menu.

```
> a:=3; b:=5; c:=22/7;
```

$$a := 3$$
$$b := 5$$
$$c := 22/7$$

Now you can use  $a$ ,  $b$  and  $c$  as though they are numbers. Thus, we have

```
> a+b;
```

$$8$$

```
> c*a;
```

$$\frac{66}{7}$$

```
> a-b^2;
```

$$-22$$

## 4.2 Variables and subsequent assignments: unassignments

If you type a variable name and a semicolon, Maple will tell you what the variable stands for:

```
> c;
```

$$\frac{22}{7}$$

```
> x;
```

$$x$$

If you specify  $y$  as an expression involving  $x$  and then assign a particular value to  $x$ , this value will get substituted into  $y$ :

```
> y:=sqrt(x+3);
```



$$y := \sqrt{x + 3}$$

```
> x:=sin(w);
```

$$x := \sin(w)$$

```
> y;
```

$$\sqrt{\sin(w) + 3}$$

If  $x$  gets changed, then  $y$  will also change. Maple remembers the original definition of  $y$ :

```
> x:=cos(u);
```

$$x := \cos u$$

```
> y;
```

$$\sqrt{\cos u + 3}$$

Notice that Maple evaluates  $x$  before telling you what the variable  $y$  stands for now. To undo this kind of thing, you can *unassign* variables. To do this for  $x$  execute:

```
> x:='x';
y;
```

$$x := x$$

$$\sqrt{x + 3}$$

We can also tell Maple to forget entirely all assignments made so far using the `restart` command.

```
> restart;
x; y; z;
```

$x$  $y$  $z$ 

Surrounding an expression in single quotes<sup>10</sup> tells Maple to delay evaluation of it. This is why `x:='x'` unassigns `x`— `x` is assigned to mean just `x` without any evaluation. Note that stopping evaluation does not stop simplification:

```
> 'x + 2 + 5';
```

 $x + 7$ 

### 4.3 Functions known to Maple. How to define your own functions

Maple has all the standard mathematical functions in its library. You can make use of them by their appropriate name. We list some of these.

- Trig and inverse trig functions: `sin`, `cos`, `tan`, `cot`, `sec`, `csc`, `arcsin`, `arccos`, `arctan`, `arccot`, `arcsec`, `arccsc`.
- Exponential and logarithm: Maple uses `exp(x)` for  $e^x$  and either `log(x)` or `ln(x)` for the natural logarithm  $\ln x$ . Logarithm in other bases are also available. For example, `log[10](x)` is  $\log_{10} x$ .
- Hyperbolic functions: `sinh`, `cosh`, etc.
- Square roots and absolute values: `sqrt(x)` and `abs(x)`.

A much more complete list can be found by asking Maple for help on the initially known functions; `?inifcns` does this.

Of course, you will find yourself in a situation where you would like to define your own functions. Suppose you are analyzing the expression  $\sqrt{1+x^2}$ , plugging several values into it and calculating the result. Maple does not have a built-in function to do this, but you can define one yourself:

---

<sup>10</sup>Be careful! Most keyboards have three different kinds of quotation marks, and they each mean something different to Maple. The single quote (') delays evaluation, the double quote (") indicates a string, and the backquote (`) indicates a name (similar to a string).

```
> f:=x->sqrt(1+x^2);
```

$$f := x \rightarrow \sqrt{1 + x^2}$$

Carefully notice the syntax. `f` is the name of the function and is written to the right side of `:=`. The function  $f$  sends the variable  $x$  to the expression  $\sqrt{1 + x^2}$ ; this is indicated by the arrow. The arrow `->` is made with the minus and greater-than keys.

After such an assignment, you can use the name `f` in the same way you use other names of functions known to Maple, e.g. `exp` or `sin`:

```
> f(0); f(1); f(-1); f(3*x+1);
```

$$\begin{aligned} & 1 \\ & \sqrt{2} \\ & \sqrt{2} \\ & \sqrt{2 + 9x^2 + 6x} \end{aligned}$$

We are so accustomed to the notation  $f(x)$  to refer to a function that we might accidentally use it in Maple. Never try to define a function with an assignment like

```
> f(x):=sqrt(x-2);
```

$$f := x \rightarrow \sqrt{x - 2}$$

It seemed like this worked, but in fact, the resulting `f(x)` does not work as expected. You will not be able to substitute into it or do anything you normally do with functions.<sup>11</sup>

Maple uses parentheses ( ) to group expressions together in a convenient form as in

```
> 2*(x+4)^3;
```

$$2(x + 4)^3$$

or to *delimit the arguments of a function*. In mathematics, we sometimes cheat and write  $\sin x$  instead of  $\sin(x)$ . In Maple, this is not permissible: you must always write `sin(x)`. Nested parentheses are permitted:

---

<sup>11</sup>Maple uses this notation for a “remember table”— you can use it to set specific values for the function. For example, if we first define  $f$  by `f:=x->sin(x)/x`, then asking for Maple for `f(0)` will give an error, since it is not a permissible value for the function. However, we can define  $f(0) = 1$  with the command `f(0):=1`; after this, `f(0)` will be 1, while  $f(x) = \sin(x)/x$  for all other values of  $x$ .

```
> f(x);
```

$$\sqrt{1+x^2}$$

```
> sin(sqrt(x+2)+ln(abs(x)));
```

$$\sin(\sqrt{x+2} + \ln|x|)$$

#### 4.4 Special constants and reserved names

Maple knows several constants. To find them out, execute the command

```
> constants;
```

$$false, \gamma, \infty, true, Catalan, FAIL, \pi$$

The names of these constants are `false`, `gamma`, `infinity`, `true`, `Catalan`, `FAIL` and `Pi`, respectively. Be careful! Maple is case sensitive. Therefore `pi` is not the same as `Pi`. For example,

```
> cos(Pi/4);
```

$$\frac{\sqrt{2}}{2}$$

```
> cos(pi/4);
```

$$\cos\left(\frac{1}{4}\pi\right)$$

Maple also uses `I` to represent the complex number  $\sqrt{-1}$ . Earlier versions used `E` for `exp(1)`; this is no longer the case. If you'd like to add `E` as a constant, you can do the following:

```
> E:=exp(1);
protect(E);
constants := constants, E;
```

$$E := e$$

$$constants := false, \gamma, \infty, true, Catalan, FAIL, \pi, e$$

While protecting your new constant `E` will stop you from accidentally redefining it, adding `E` to the list of known constants doesn't really affect things very much.

There are many commands and special objects in Maple that have preassigned names. In creating your own expressions and setting the name of your own variables, you cannot use any of the following 30 keywords:

```
and    by    do      done   elif   else   end    fi     for    from
if     in     intersect local  minus  mod    not    od     option options
or     proc  quit     read   save   stop   then   to     union  while
```

You can check this list for yourself using the help system by giving the command `?keywords`.

Maple also starts with a large number of predefined functions and constants. Many of these can be listed with `inifcns` and `ininames`. If you try to use one by mistake, you will receive an error:

```
> gamma:=x^2+5;
```

```
Error, attempting to assign to name 'gamma' which is protected
```

If you insist, you can use `unprotect` to allow you to assign to this name. But **don't do this** unless you know what you are doing; this changes the meaning of this name even when it is used internally, and can have unexpected consequences.

## 5 The limit command

The limit command is used to compute limits. Its syntax is basically

```
> limit(f,x=a);
```

where `f` is a Maple algebraic expression which (in general) depends on the variable `x` and `a` is the expression which `x` approaches. For example, to compute

$$\lim_{x \rightarrow 2} \frac{3x - 6}{x^2 - 4},$$

we execute the Maple command

```
> limit((3*x-6)/(x^2-4),x=2);
```

$$\frac{3}{4}$$

Maple can also deal with limits which do not exist, for example  $\lim_{x \rightarrow 0} \frac{1}{x}$ :

```
> limit(1/x,x=0);
```

*undefined*

In the following exercise, we define the slope of a straight line passing through the points  $(x_1, y_1)$ ,  $(x_2, y_2)$ , use this function to find the slope  $m(x)$  of the line passing through the points  $(1, 1)$  and  $(x, x^2)$ , and finally compute the limit of  $m(x)$  as  $x \rightarrow 1$ :

```
> slope:=(x1,y1,x2,y2)->(y2-y1)/(x2-x1);
  m:=x->slope(1,1,x,x^2);
```

$$\text{slope} := (x_1, y_1, x_2, y_2) \rightarrow \frac{y_2 - y_1}{x_2 - x_1}$$

$$m := x \rightarrow \text{slope}(1, 1, x, x^2)$$

```
> limit(m(x),x=1);
```

2

Maple is sometimes unable to determine a limit or whether it exists. In such a case, it will return nothing after you execute the limit command.

Maple computes *two-sided* limits. For example, if you specify the argument  $x=1$  as in the previous example, Maple assumes you mean that  $x$  approaches 1 from either the right or the left through real values only (as opposed to complex ones). However, Maple can compute one sided and complex limits also:

```
> limit(1/x,x=0,right);
```

$\infty$

```
> limit(1/x,x=0,left);
```

$-\infty$

```
> limit(x*log(x),x=0,complex);
```

0

You can also compute limits as  $x \rightarrow \infty$  and limits of functions of more than one variable:

```
> limit(arctan(x),x=infinity);
```

$$\frac{\pi}{2}$$

```
> limit(x/(x^2+y^2),{x=0,y=0});
```

*undefined*

A common error is to try to compute the limit of a function  $f(x)$  when  $x$  has been previously given a value. If you find yourself in this situation, unassign the value of  $x$  executing the command  $x:=x'$ ; . The previous example will fail if  $y$  still has the value assigned to it in section 4.2.

## 6 The diff (and Diff) command

The `diff` command is used to compute derivatives of Maple expressions. Its syntax is basically

```
> diff(f,x);
```

where  $f$  is an algebraic expression and  $x$  is the variable with respect to which the derivative is taken. For example, to compute

$$\frac{d}{dx} \frac{3x-6}{x^2-4},$$

we execute the Maple command

```
> diff((3*x-6)/(x^2-4),x);
```

$$\frac{3}{x^2-4} - 2 \frac{(3x-6)x}{(x^2-4)^2}$$

You must specify the variable with respect to which you wish to take the derivative because, almost all the time, there are constants and other parameters around:

```
> diff(exp(a*x),x);
```

$$ae^{ax}$$

Maple knows theoretical results about derivatives. For example, even if the functions  $g(x)$  and  $h(x)$  are not defined, Maple will give you the derivative of their product in terms of the functions and their derivatives:

```
> diff(h(x)*g(x),x);
```

$$\left(\frac{\partial}{\partial x}h(x)\right)g(x) + h(x)\left(\frac{\partial}{\partial x}g(x)\right)$$

## 6.1 Higher order derivatives

To compute higher order derivatives, you can of course iterate the above command:

```
> diff(diff(3*sin(x),x),x);
```

$$-3 \sin(x)$$

Such a command first calculates the derivative with respect to  $x$  of the expression `3*sin(x)` and then differentiates the result to obtain the second derivative. The same result can be accomplished with either one of the following commands:

```
> diff(3*sin(x),x,x);
```

$$-3 \sin(x)$$

```
> diff(3*sin(x),x$2);
```

$$-3 \sin(x)$$

There is an obvious extension of this command when calculating more than two derivatives. You can also calculate partial derivatives of higher order:

```
> diff(ln(x)*exp(-2*x),x$3);
```

$$2\frac{e^{-2x}}{x^3} + 6\frac{e^{-2x}}{x^2} + 12\frac{e^{-2x}}{x} - 8\ln(x)e^{-2x}$$

```
> diff(exp(x)*y^2*sin(z),x,y,z);
```

$$2e^x y \cos(z)$$



## 6.2 Implicit differentiation

Maple knows how to take derivative of both sides of an equation, and, as indicated above, Maple can also take derivatives “theoretically” even if the definition of the functions involved are not specified. However, you must always indicate *explicitly* the dependence of the functions on their variables. So, if you are thinking of  $y$  as a function of  $x$ , you must write  $y(x)$  in your equation rather than just  $y$ .

Here is a typical implicit differentiation problem. Consider the equation  $x^2y - 3y^3x = 0$ . Find the slope of the graph of the curve defined by this equation at the point  $(3, 1)$ :

```
> eq:=x^2*y(x)-3*y(x)^3*x = 0;
   deq:=diff(eq,x);
   solve(deq,diff(y(x),x));
```

$$\begin{aligned} eq &:= x^2y(x) - 3y(x)^3x = 0; \\ deq &:= 2xy(x) + x^2 \left( \frac{\partial}{\partial x} y(x) \right) - 9y(x)^2 \left( \frac{\partial}{\partial x} y(x) \right) - 3y(x)^3 = 0 \\ &\quad \frac{2xy(x) - 3y(x)^3}{x^2 - 9y(x)^2x} \end{aligned}$$

```
> subs(y(x)=1, x=3, %);
```

$$\frac{1}{6}$$

The first command defines the equation relating  $x$  and  $y$ , the second defines an equation `deq` giving the relation between  $x$ ,  $y(x)$  and the derivative of  $y$  with respect to  $x$ . Using `solve`, we solve `deq` to find the derivative as a function of  $x$  and  $y(x)$ , and in the result we substitute  $x=3$  and  $y=1$  to obtain the desired slope (see section 8 to learn about the `subs` command).<sup>12</sup>

Maple actually has a built-in command to do implicit differentiation, `implicitdiff`. Thus, we could have done this same problem more concisely using the single command

```
> subs({y=1,x=3}, implicitdiff(x^2y-3y^3x=0,y,x));
```

$$\frac{1}{6}$$

Occasionally, to make your worksheets easier to read, you may wish to have Maple display a derivative in standard mathematical notation without evaluating it. For this purpose there

---

<sup>12</sup>Again, the major caveat with the Maple derivative command is to avoid calculating derivatives with respect to a variable that has been previously given a value. In such an instance, unassign the value of the variable first.

is an *inert* capitalized form of the `diff` command: `Diff`. The two forms `diff` and `Diff` are usually combined to produce meaningful sentences:

```
> Diff(exp(x)/(1-x),x);
```

$$\frac{\partial}{\partial x} \frac{e^x}{1-x}$$

```
> Diff(exp(x)/(1-x),x)=diff(exp(x)/(1-x),x);
```

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

This command is particularly useful to produce easy to read worksheets that involve partial derivatives:

```
> f:=x^3*exp(y)-sin(x*y);
```

$$f := x^3 e^y - \sin(xy)$$

```
> Diff(f,x,x,y)=diff(f,x,x,y);
```

$$\frac{\partial^3}{\partial y \partial x^2} x^3 e^y - \sin(xy) = 6x e^y + \cos(xy) x y^2 + 2 \sin(xy) y$$

## 7 The int (and Int) command

### 7.1 Symbolic integrals

The `int` command is used to compute both definite and indefinite integrals of Maple expressions. Its syntax is basically

```
> int(f,x);
```

where `f` is an algebraic expression and `x` is the integration variable. For example, to compute

$$\int \frac{3x-6}{x^2-4} dx,$$

we execute the Maple command

```
> int((3*x-6)/(x^2-4),x);
```

$$3 \ln(x + 2)$$

Notice that Maple does not provide the constant of integration. You will occasionally have to take this into account and provide your own constant.

You must specify the variable of integration. In expressions involving other parameters, Maple assumes that you want the integral of the expression as the variable you specify changes and that all other parameters in the expression represent constants:

```
> int(exp(a*x),x);
```

$$\frac{e^{ax}}{a}$$

To compute a definite integral, the *range* over which the integration variable moves must be specified:

```
> int(x^2*exp(x),x=0..2);
```

$$2e^2 - 2$$

The `int` command in Maple has a particular behavior in certain situations:

- If Maple cannot compute the integral in *closed form*, it will return it unevaluated:

```
> int(ln(sin(sqrt(x^12-5*x^7+50*x+2))),x);
```

$$\int \ln(\sin(\sqrt{x^{12} - 5x^7 + 50x + 2})) dx$$

- Sometimes the integral cannot be evaluated in closed form in terms of *elementary* functions, but the answer has a special name in mathematical circles due to its importance in applications. For example,

```
> int(sin(2*x)/x,x);
```

$$Si(2x)$$

Here `Si` is the *special name* of one of these functions which appear frequently in mathematical physics. To learn more about it, ask Maple for help on it with the command `?Si`. This will bring up a window with information about the function `Si`. If Maple responds

to an integral with one of these functions, it is quite likely that the integral cannot be evaluated in terms of elementary functions.

- Sometimes, a result may be expressed in terms of the roots of a polynomial which does not factor over the rationals. For example,

```
> int(1/(x^8 + 1), x);
```

$$\sum_{-R=\text{RootOf}(1+16777216-Z^8)} -R \ln(x + 8-R)$$

In this answer, the sum is taken over all roots  $R$  of the polynomial  $1 + 16777216z^8$ , and the summand is  $R \log(x + 8R)$ .

As with `diff`, there is an *inert* form `Int` of the integral command which can be used in combination with `int` to produce easily readable worksheets:

```
> Int(ln(1+3*x), x=1..4);
```

$$\int_1^4 \ln(1 + 3x) dx$$

```
> Int(ln(1+3*x), x=1..4)=int(ln(1+3*x), x=1..4);
```

$$\int_1^4 \ln(1 + 3x) dx = \frac{13}{3} \ln(13) - 3 - \frac{4}{3} \ln(4)$$

The inert form `Int` is also very useful in many situations when you wish to delay the evaluation of an integral, as we shall see below.

## 7.2 Numerical integration

With the `evalf` command, you can force Maple to apply a numerical approximation technique for definite integration:

```
> evalf(Int(sqrt(1+x^10), x=0..1));
```

1.040899075

Notice that the we are using the inert form of the integration command, `Int`. This prevents Maple from attempting to evaluate the integral symbolically and then applying `evalf` to the answer. In many cases, this can save a huge amount of time, because Maple will work very hard

to try to compute the symbolic form of the integral. For example, approximating  $\int_0^1 e^{\sin(x)} dx$  with the command

```
> evalf(int(exp(sin(x)),x=0..1);
```

took more than 10 times the amount of time needed to execute

```
> evalf(Int(exp(sin(x)),x=0..1);
```

More complicated integrals can have even more dramatic differences.

### 7.3 Approximations through Riemann Sums

Maple can also compute expressions that approximate a definite integral using rectangles, trapezoids, etc. It can also do approximations using Simpson's rule. The only difficulty is that these commands are contained in a separate library called `student` and you must load this library into computer memory before you can use it:

```
> with(student):
```

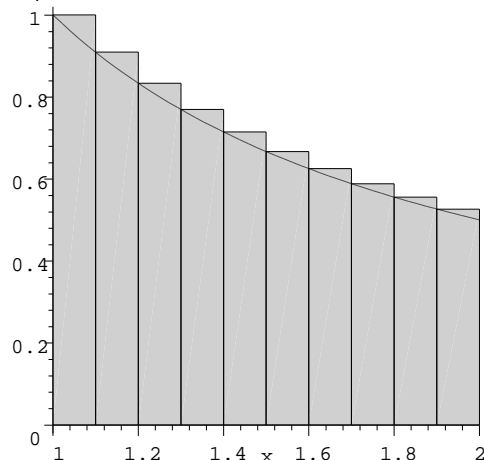
Once loaded, you can play with it. For example, if you want to use the left sum approximation to an integral, the height of each rectangle is determined by the value of the function at the left side of each interval. You may specify the number of intervals you wish to use. If you do not, Maple will use four intervals by default:

```
> leftsum(x^4*ln(x),x=1..4,10);
```

$$\frac{3}{10} \left( \sum_{i=0}^9 \left(1 + \frac{3}{10}i\right)^4 \ln\left(1 + \frac{3}{10}i\right) \right)$$

Maple can also draw pictures of the rectangles used to approximate the integral of a function on a given interval:

```
> leftbox(1/x,x=1..2,10);
```



The following is the result of using the Simpson's rule for the function  $x\sin(x^2)$  on the interval  $1 \leq x \leq 5$ . Since the number of intervals is not specified, Maple assumes four equal intervals by default:

```
> simpson(x*sin(x^2),x=1..5);
```

$$\frac{1}{3} \sin(1) + \frac{5}{3} \sin(25) + \frac{4}{3} \left( \sum_{i=1}^2 2i \sin(4i^2) \right) + \frac{2}{3} \left( \sum_{i=1}^1 (1 + 2i) \sin((1 + 2i^2)^2) \right)$$

## 7.4 Multiple integrals

Maple does iterated integrals:

```
> int(int(x^2*y^3,x=0..y),y=2..3);
```

$$\frac{2059}{21}$$

This is literally an iterated integral: the integral of  $x^2y^3$  with respect to  $x$  on the interval  $[0, y]$  is nested inside the integral with respect to  $y$  on the interval  $[2, 3]$ , in effect calculating

$$\int_2^3 \int_0^y x^2 y^3 dx dy.$$

## 8 The subs command

You have learned in the previous sections that if you assign a particular value to a variable, Maple will remember it for the rest of the session until you unassign it. At times, you might like to *plug in* some values into an expression without altering the variable (or variables) involved. Maple has the `subs` command which reports the result of making such a substitution:

```
> y:=x+3;
```

$$y := x + 3$$

```
> subs(x=2,y);
```

Notice that Maple reports the value obtained when plugging in  $x=2$  into  $y$  without giving that value to  $x$  or altering  $y$  in any way:

```
> x;
```

$$x$$

```
> y;
```

$$x + 3$$

You can make more than one substitution:

```
> z:=x+sin(2*h):
  subs(x=3*u+1,h=4*Pi*q,z);
```

$$3u + 1 + \sin(8\pi q)$$

In this example, the first assignment is ended with a colon `:` to suppress the display. You can do this whenever you want to make an assignment which you do not want to see displayed. Then, plugging  $x = 3u + 1$  and  $h = 4\pi q$  into the expression for  $z$ , we obtain  $3u + 1 + \sin(8\pi q)$ . Notice that the values of `x` and `h` that we plug in do not need to be enclosed in braces when issuing the command.

You should be aware of a few things when using the `subs` command:

- You need not fear *circular substitutions*:

```
> z;
```

$$x + \sin(2h)$$

```
> subs(x=2*x,z);
```

$$2x + \sin(2h)$$

On the other hand, long chains of self-referential substitutions may produce results which are hard to predict:

```
> subs(h=x,x=z,z);
```

$$x + \sin(2h) + \sin(2x + 2\sin(2h))$$

Simultaneous substitutions are made enclosing the values you plug in in braces. This is in contrast with the left-to-right substitution above:

```
> subs({h=x, x=z}, z);
```

$$x + \sin(2h) + \sin(2x)$$

- In certain situations, `subs` will not be able to find the expression you are substituting for if it is complicated. In those cases, try to *reword* your substitution request so that Maple *understands it*, or break the task into several manageable ones. In other cases, you might consider using `algsubs`, which allows you to substitute one algebraic expression for another.
- `subs` does what is called “syntactic substitution”; it substitutes one set of symbols for another. If you are using it to evaluate an expression with certain values “plugged in”, using `eval` might be more appropriate. Here is an example that should illustrate the distinction:

```
> integral := int(f(t,a), t=a..x);
```

$$\int_a^x f(t, a) dt$$

```
> eval(integral, {t=0, a=1});
```

$$\int_1^x f(t, 1) dt$$

```
> subs({t=0, a=1}, integral);
```

$$\int f(0, 1) d0 = 1..x$$

The seemingly nonsensical answer using `subs` above is because Maple interpreted the expression as meaning `int(f(0,1), 0=1..x)` — what you get by replacing `t` by 0 and `a` by 1 without trying to interpret the expression.

## 8.1 Why is `subs` useful?

Often it is better to use `subs` rather than to make assignments to the variables in an expression. For example, if



```
> y:=3*x^2+5*x+2;
```

$$y := 3x^2 + 5x + 2$$

and you want to compute the difference quotient  $(y(x+h) - y(x))/h$ , then you should execute

```
> dq:=(subs(x=x+h,y)-y)/h;
```

$$dq := \frac{3(x+h)^2 + 5h - 3x^2}{h}$$

because if we set

```
> x:=x+h;
```

**Warning: Recursive definition of name**

$$x := x + h$$

any call to `y` will result in

```
> y;
```

**Error, too many levels of recursion**

In earlier versions of Maple, this action would have crashed your session.

On the other hand,

```
> x:=4;
```

$$x := 4$$

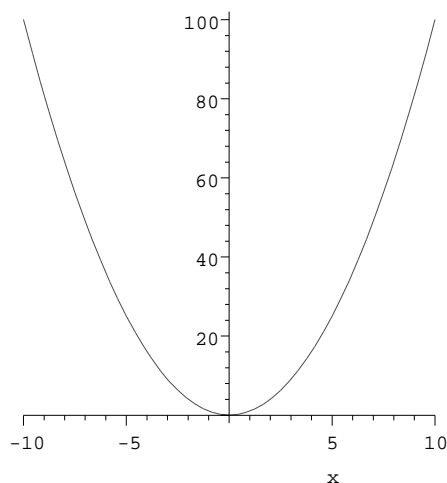
```
> y;
```

is alright, but now `y` is a number and we cannot plug any other values of `x` into it. Using `subs` is a better way to do the plugging in, because we need not unassign `y`.

## 9 Plotting with Maple

Let us now use the graphics facilities of Maple:

```
> plot(x^2, x=-10..10);
```



Maple displays the graph of the function  $x^2$  from  $x = -10$  to  $x = 10$  as requested. In the worksheet interface, clicking the right mouse button on the graph pops up a menu with various options which allow you to manipulate your graph, and/or change its appearance in a suitable way. Clicking the left button “selects” the graphic, drawing a black box with “handles” around it. You can drag the mouse button on these to resize the plot. Also, the coordinates of the point where you clicked appear in the upper left of the worksheet. Selecting a graphic also changes the context bar: shortcuts for manipulating the axes, aspect ratio, and plotting style appear as buttons along the top.

The basic syntax of the plot command is

```
> plot(f, range);
```

where **f** is the expression to be plotted and **range** is the range of the parameter(s) for which you would like to see the plot of **f**. In the example above, we indicated the range of **x** by **x=-10..10**. Maple automatically chooses a scale on the vertical axis. There are many options to the plot command; while we will mention a few of them, the on-line help system has all the details.

Try both of the two similar plot commands below:

```
> plot(x^2-x, x=-1..2, y=-1..2);
plot(x^2-x, x=-1..2, -1..2);
```

In both cases, the plot is displayed for values of  $y$  between  $-1$  and  $2$ , but in the first, the variable  $y$  is explicitly named in the command, and the vertical axis gets a label. Whenever you are plotting an expression, the variable in the domain of the function *must* be specified.

## 9.1 Plotting several functions (or curves) on the same axes

We can plot more than one function (or curve) on the same coordinate system. Consider the following example, recalling the slope function of the line passing through the points  $(x_1, y_1)$ ,  $(x_2, y_2)$ , used to define the slope  $m(x)$  of the line passing through  $(1, 1)$  and  $(x, x^2)$ :

```
> slope:=(x1,y1,x2,y2)->(y2-y1)/(x2-x1);
m:=x->slope(1,1,x,x^2);
```

$$\text{slope} := (x1, y1, x2, y2) \rightarrow \frac{y2 - y1}{x2 - x1}$$

$$m := x \rightarrow \text{slope}(1, 1, x, x^2)$$

We now define the equation of the lines passing through  $(1, 1)$  and  $(x, x^2)$  for  $x = 2, 3, 4, 5$  and plot the four lines together with the function  $x^2$ .

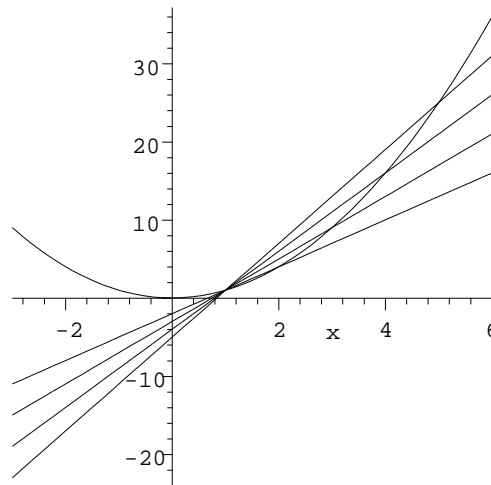
```
> line1:=m(2)*(x-1)+1;
line2:=m(3)*(x-1)+1;
line3:=m(4)*(x-1)+1;
line4:=m(5)*(x-1)+1;
plot({x^2,line1,line2,line3,line4},x=-3..6);
```

$$\text{line1} := 3x - 2$$

$$\text{line2} := 4x - 3$$

$$\text{line3} := 5x - 4$$

$$\text{line4} := 6x - 5$$



Observe that the expressions to be plotted are *enclosed* in braces  $\{ \}$ . Besides specifying the range on which  $x$  varies as done above, you can also specify and label the  $y$  range proceeding exactly in the same manner as when plotting a single function.<sup>13</sup>

---

<sup>13</sup>We could have done this all in a single plot command, either explicitly typing each of the lines, or using `seq` to generate them: `plot({x^2, seq( m(i)*(x-1)+1, i=2..5)}, x=-3..6);`

## 9.2 Fancier plotting

The `plot` command is very powerful, and you should look for details about it using the online help facility; the command `?plot` will bring this up. There are also several related pages, and some instructive example worksheets that are well worth looking at.

To make you aware of some of its features, we will discuss a few of them here:

- It is possible to have Maple *plot points*. This can be done in two ways depending upon how the points are generated. If you have a list of specific points, you can assign them to a name and then plot them as follows:

```
> points:=[[1,2],[1.5,1],[2,-1],[2.5,0.5],[3,1],[3.5,0.6],[4,0.2]]:
   plot(points,style=POINT);
```

If `style=POINT` is replaced by `style=LINE`, the dots will be connected by segments of lines. In the worksheet interface, you can switch from one style to the other by first clicking the left mouse button on the graph, and then hitting the relevant button on the buttonbar.

On the other hand, if the points result from evaluating an expression at several values of  $x$ , you can use `plot` in its usual form, but specifying `style=POINT`:

```
> plot(x^2,x=-2..2,style=POINT);
```

Approximately 50 points will be plotted this way. You may notice that they will not usually be equally spaced, but concentrated in areas where the graph curves more. You can insist that more points be plotted by using the `numpoints` option:

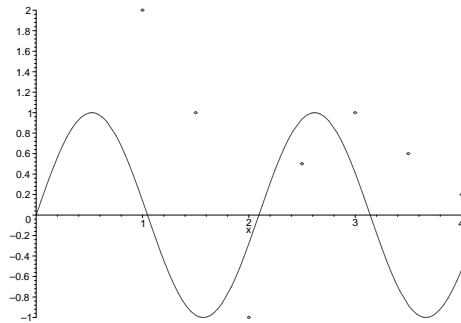
```
> plot(x^2,x=-2..2,style=POINT,numpoints=150);
```

- Several more specialized plotting functions are contained in a separate library called `plots`. For basic plotting, there are two commands from this library that are especially useful. They are `textplot` and `display`. To load these commands into the computer memory, execute the statement

```
> with(plots,textplot,display);
```

The `display` command is useful to combine different kinds of plots into one picture. You can merge standard plots of expressions, plots of points, plots of text (This is the utility of `textplot`: to label things), and animations. In the next example, we combine the point plot of the variable `points` defined above with a plot of `sin(3*x)`. We first define them as different plots and then display them together. Note that except for the third command, we end with a colon (`:`) in order to suppress the Maple output.

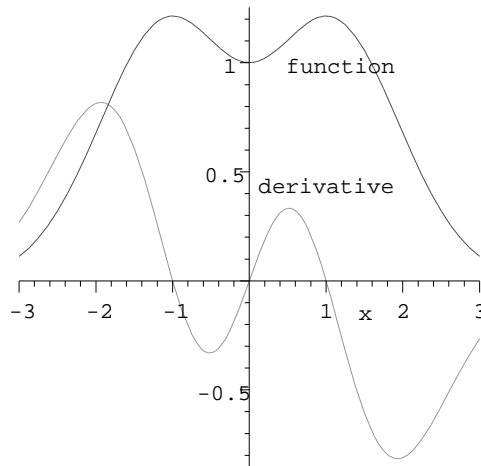
```
> plot1:=plot([points],style=POINT):
plot2:=plot(sin(3*x),x=0..4):
display({plot1,plot2},view=[0..4,-1..2]);
```



The `view` option controls the horizontal and vertical ranges to be displayed.

You can attach labels to your plots by combining `textplot` with `display`. The basic syntax of `textplot` is to use it with an argument of the form `[a,b,'name']` which places the word `name` inside the quotes on the plot so that it is centered at the point  $(a,b)$ . For example:

```
> y:=(1+x^2)*exp(-x^2/2): d:=diff(y,x):
F:=plot({y,d},x=-3..3):
G:=textplot({[1,1,'function'],[0.75,0.45,'derivative']}):
display({F,G});
```



We end with a word about plotting functions (as opposed to expressions), and a situation in which it is a good idea to do so. Sometimes you will have a relationship that you want to plot in the form of a function rather than an expression:

```
> f:=x->x^3*exp(-x):
```

In this situation, you can plot  $f(x)$  as explained above. Alternatively, you may use:

```
> plot(f, 0..3);
```

The two ways of plotting a function should not be confused. Neither of the following two statements will work correctly:

```
> plot(f(x), 0..3);
```

Plotting error, empty plot

```
> plot(f, x=0..3);
```

As you see, the first gives an error; the second gives a plot with nothing plotted on it.

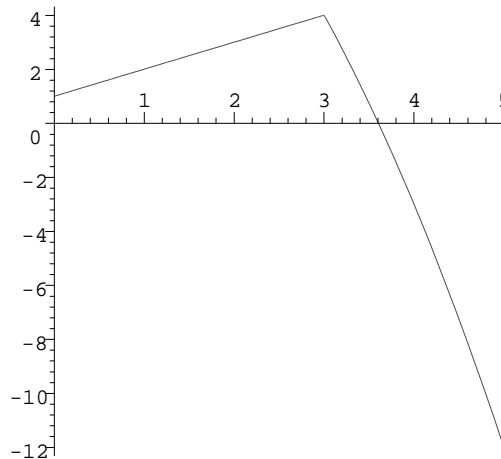
If you have a function defined using an *if-then* clause, you must use the function plotting command:

```
> f:=x-> if x<3 then x+1 else -x^2+13 fi:
  plot(f(x),x=0..5);
```

Error, (in f) cannot evaluate boolean

The error resulted because Maple attempts to understand the expression  $f(x)$  before it has a value for  $x$ . One way around this<sup>14</sup> is to use the command

```
> plot(f, 0..5);
```



You can see from the graph that the function  $f$  is probably continuous but not differentiable at  $x=3$ .

## 10 Exercises

1. Evaluate  $\pi^{\sqrt{2}}$  to 30 decimal places.
2. Input the following as a Maple function:

$$f(x) = \frac{\sqrt{x}(-116x^3 + 8x^4 + 558x^2 - 891x)}{2x^5 - 29x^4 + 140x^3 - 225x^2}.$$

---

<sup>14</sup>We could also delay evaluation of  $f(x)$  using quotes, with `plot('f(x)', x=0..5);`. Alternatively, we could define  $f$  using `piecewise`, as `f:=x-> piecewise(x<3, x+1, -x^2+13);`.

- a) Write this in a simpler form (that is, factor and reduce it).
  - b) Draw the graph of the function for  $x$  between 0 and 1). Adjust the vertical range so that some detail can be seen.
  - c) Compute the area of the part of the curve that lies above the  $x$ -axis and between  $x = 0$  and  $x = 5$  (that is, integrate the function over the appropriate range of  $x$  values). Give your answer both in actual form and as a decimal approximation to about 10 places.
  - d) What is the value of the integral if you remove the factor of  $\sqrt{x}$  from the numerator?
  - e) Use the derivative of  $f(x)$  to determine for what real value  $x$  the function has a local maximum. Use an approximation to about 8 decimal places.
3. Use the commands `seq` and `ithprime` to generate a list of the first 20 primes. Compute the sum of these 20 primes, and give its integer factors.
  4. Find the solutions of the system of equations  $\{x^2 - y^2 = 4, x - 2y = 2\}$ .
  5. Draw a graph showing both  $\cos x$  and its fifth Taylor polynomial (that is,  $1 - \frac{1}{2!}x^2 + \frac{1}{4!}x^4$ ) for  $x$  between  $-4$  and  $4$ . How many terms do you seem to need to get *good agreement* in this range.

**Hint:** use a variation of the command `convert(taylor(cos(x),x,5),polynom)` to make this work.

Think of a suitable way to demonstrate that the approximation you have taken is *good*.

# If the Curve Fits, Wear It

It is not uncommon in science and mathematics to have a collection of data points

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

for which we believe there is a functional relationship between the  $x_i$  and the  $y_i$ . Often, we want to get some idea of what might be happening at points other than those for which we have data. That is, we want to find some  $f$  so that  $f(x_i) = y_i$  for each of our points. Of course, in order to do this, we will need to have some knowledge (or make some guesses) about the function  $f$ , or sometimes about how precisely the points are known. In this chapter, we'll look at several different approaches to this problem of fitting a curve through (or near) some given data.

---

## Subsections

- [Interpolation](#)
  - [Polynomial Interpolation](#)
  - [Connect-the-Dots and Splines](#)
- [When the data is approximate](#)
- [Fitting a line to data](#)
- [Fitting a cubic to data](#)
- [Fitting other types of functions](#)
- [Fitting a circle](#)
- [Robust fitting](#)
- [A nod toward statistics](#)



# Chapter 2

## If the Curve Fits, Wear It

It is not uncommon in science and mathematics to have a collection of data points

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

for which we believe there is a functional relationship between the  $x_i$  and the  $y_i$ . Often, we want to get some idea of what might be happening at points other than those for which we have data. That is, we want to find some  $f$  so that  $f(x_i) = y_i$  for each of our points. Of course, in order to do this, we will need to have some knowledge (or make some guesses) about the function  $f$ , or sometimes about how precisely the points are known. In this chapter, we'll look at several different approaches to this problem of fitting a curve through (or near) some given data.

### 1 Interpolation

First, we will look at the situation where we the data points are assumed to be known exactly (or at least to good enough precision that we ignore any errors), and we want to find a curve of some type is chosen that passes through each of the data points. This practice is called *interpolation*. In the latter part of this chapter (beginning with §2), we will relax the restriction that the curve pass exactly through the points.

Of course, there are serious restrictions involved in interpolation. Since a line is determined by two points, if we have more points, we must either use a higher degree polynomial, or use several line segments. If we have  $n$  data points we wish to interpolate, either we can fit a polynomial of degree  $n - 1$  or find  $n - 1$  line segments which “connect the dots”. We'll address both of these in this section, as well as a method which is a cross between the two approaches.

#### 1.1 Polynomial Interpolation

Given  $n$  data points  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ , there is a unique polynomial of degree  $n - 1$  passing through them. Finding this polynomial is just a matter of solving the  $n$  linear equations

for the coefficients. For example, suppose we have 4 data points  $(x_1, y_1)$ ,  $(x_2, y_2)$ ,  $(x_3, y_3)$ , and  $(x_4, y_4)$ . There is exactly one degree 3 polynomial  $p(x) = ax^3 + bx^2 + cx + d$  which passes through all four points. Plugging in our data yields the equations

$$\begin{aligned} ax_1^3 + bx_1^2 + cx_1 + d &= y_1 \\ ax_2^3 + bx_2^2 + cx_2 + d &= y_2 \\ ax_3^3 + bx_3^2 + cx_3 + d &= y_3 \\ ax_4^3 + bx_4^2 + cx_4 + d &= y_4 \end{aligned}$$

These are easily solved for the coefficients  $a$ ,  $b$ ,  $c$ , and  $d$ .

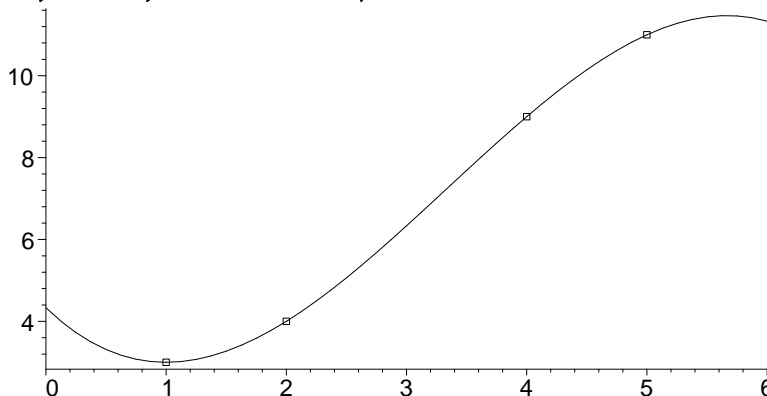
Maple 7<sup>1</sup> has a built-in command called `PolynomialInterpolation` to do this automatically for us, as part of the `CurveFitting` package. We'll do a brief example.

```
> data:=[[1,3],[2,4],[4,9],[5,11]];

with(CurveFitting):
cubfit:=PolynomialInterpolation(data,x);
```

$$cubfit := -\frac{1}{6}x^3 + \frac{5}{3}x^2 - \frac{17}{6}x + \frac{13}{3}$$

```
> plots[display](
  plot(data,style=point,color=black,symbol=BOX),
  plot(cubfit,x=0..6,thickness=2));
```



That seems to have worked pretty well. However, we should remark that polynomials are somewhat inflexible: minor variations in the data can have drastic effects on the behavior between the points. To emphasize this, we'll generate 10 points along the line  $y = 2x + 1$ , with one more point in the middle that is a bit above the line.

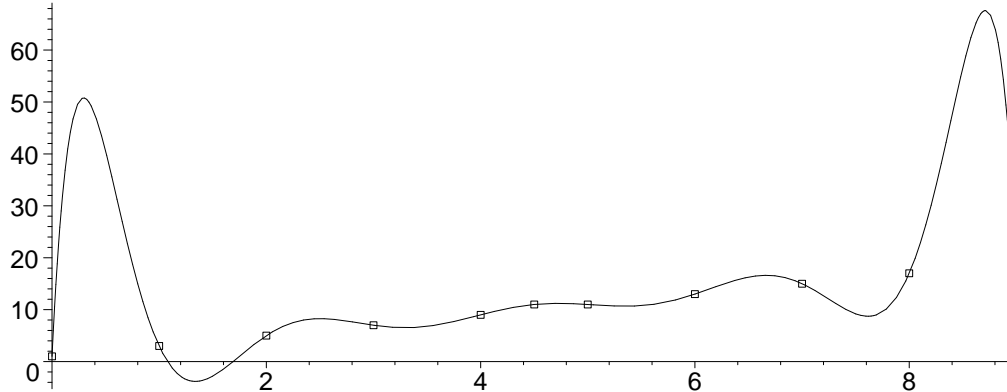
```
> data2:= [seq([i,2*i+1],i=0..4),[9/2,11],seq([i,2*i+1],i=5..9)];
```

---

<sup>1</sup>earlier versions of Maple have a similar command called `interp` with slightly different syntax.

```
data2 := [[0, 1], [1, 3], [2, 5], [3, 7], [4, 9], [9/2, 11], [5, 11], [6, 13], [7, 15], [8, 17], [9, 19]]
```

```
> plots[display](
  plot(data2, style=point, color=black, symbol=BOX),
  plot(PolynomialInterpolation(data2, x), x=0..9, thickness=2));
```



This doesn't look much like a straight line, does it? Even though there is only one data point which is 1 unit above the line  $y = 2x + 1$ , at  $x = 1/2$  the polynomial is about 50 units above this line.

## 1.2 Connect-the-Dots and Splines

As we said at the start of this section, another choice is not to use a single function, but to use several segments which connect the dots. The result would be a piecewise-affine function consisting of several line segments defined on a number of intervals. In the example using `data2` from the previous section, this function would be

$$f(x) = \begin{cases} 2x + 1 & \text{if } x \leq 4 \\ 4x - 7 & \text{if } 4 < x \leq \frac{9}{2} \\ 11 & \text{if } \frac{9}{2} < x \leq 6 \\ 2x + 1 & \text{if } 6 < x \end{cases}$$

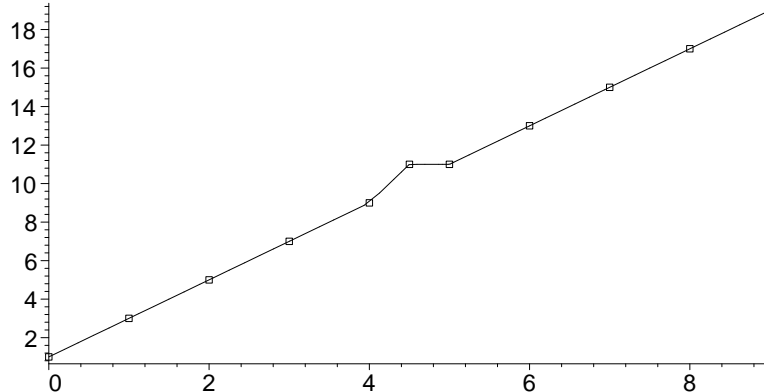
Although writing a maple procedure to figure this out for us would not be hard, such a piecewise affine curve is a version of something called a spline, which we will discuss more soon. This type of connect-the-dots curve is a linear (or degree 1) spline, and Maple can find it using the `Spline` command out of the `CurveFitting`<sup>2</sup> package.

```
> lspline:=Spline(data2,x,degree=1):
```

---

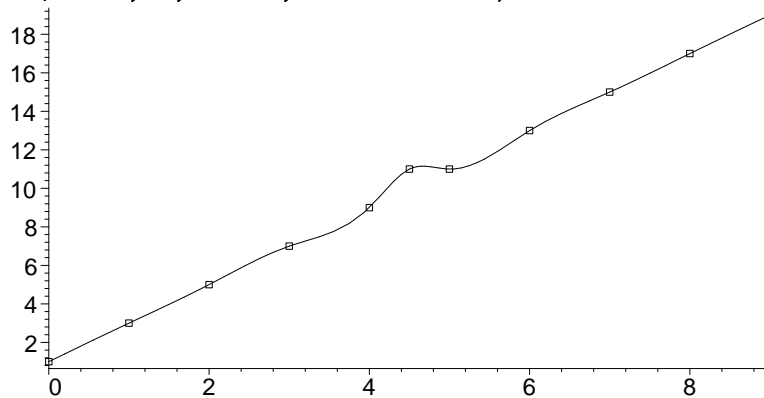
<sup>2</sup>versions prior to Maple 7 have a similar command called `spline` which takes slightly different arguments.

```
> plots[display](
  plot(data2,style=point,color=black,symbol=BOX),
  plot(lsplines,x=0..9,thickness=2));
```



If we leave off the `degree=1`, Maple gives us a cubic spline instead:

```
> plots[display](
  plot(data2,style=point,color=black,symbol=BOX),
  plot(Spline(data2,x),x=0..9,thickness=2));
```



Notice how the cubic spline fits the points about as well as the “connect-the-dots” piecewise affine curve, but has no corners. What is this curve that Maple is showing us?

To make a spline, instead of connecting each pair of points by a straight segment, we put in a polynomial of degree  $d$ . We must ensure that the polynomial passes through the two endpoints of each segment, but this still leaves us an additional  $d - 1$  conditions per polynomial segment. The entire curve can be made smoother by choosing each polynomial segment so that the first  $d - 1$  derivatives at one endpoint agree with the derivatives of the previous segment. Since there is no segment before the first or after the last, this leaves  $d - 1$  more conditions still to specify. For a “natural spline”, the high order derivatives at the endpoints are set to zero:

$$p^{(d-1)}(x_1) = p^{(d-2)}(x_1) = \dots = p^{(\frac{d-1}{2})}(x_1) = 0$$

and

$$p^{(d-1)}(x_n) = p^{(d-2)}(x_n) = \dots = p^{(\frac{d-1}{2})}(x_n) = 0$$

(if  $d - 1$  is odd, we have one fewer zero derivative at  $x_n$  than at  $x_1$ ). Thus, for a cubic spline, the second derivatives vanish at eachpoint, and the polynomial segments have the same value, slope, and concavity at each interior point.

You might think that using higher degree polynomials in between would give a better appearance, but the same issues that arose in earlier begin to show up. Note that if we have  $n$  data points, fitting a spline of degree  $n - 1$  is exactly the same as the interpolating polynomial of degree  $n - 1$ . Cubic splines are the most widely used, because they look quite smooth to the eye, but have enough freedom that they don't have to wiggle too much in order to pass through each data point. Many computer drawing programs have the ability to fit cubic splines through a set of points.

## 2 When the data is approximate

Now we relax the restriction that the function we are searching for must pass exactly through each of the data points. This is the typical situation in science, where we make a measurement or do an experiment to gather our  $y$  values from the input  $x$  values.

More specifically, let's assume that  $y = f_{c_1, \dots, c_k}(x_1, \dots, x_m)$  is a real-valued function of  $(x_1, \dots, x_m)$  which depends upon  $k$  parameters  $c_1, \dots, c_k$ . These parameters are unknown to us. However, suppose that we can perform repeated experiments that for given values of  $(x_1, \dots, x_m)$  allows us to measure output values for  $y$ . How can we estimate the parameters  $c_1, \dots, c_k$  that best correspond with this information?

Let us assume that the experiment measuring the value  $y = f_c(x)$  for specific input values  $x = (x_1, \dots, x_m)$  is repeated  $n$  times. We will then obtain a system of equations

$$\begin{aligned} f_{c_1, \dots, c_k}(x_{11}, x_{12}, \dots, x_{1m}) &= y_1 \\ f_{c_1, \dots, c_k}(x_{21}, x_{22}, \dots, x_{2m}) &= y_2 \\ &\vdots \\ f_{c_1, \dots, c_k}(x_{n1}, x_{n2}, \dots, x_{nm}) &= y_n \end{aligned}$$

where  $x_{ij}$  and  $y_i$  are the measurements of  $x_j$  and  $y$  in the  $i^{\text{th}}$  experiment. These experimental results gives us information about the unknown coefficients  $c_i$ .

Since we may perform the experiments as many times as we wish, we may end up with more relations of the type above than unknowns (that is,  $n$  is larger than  $k$ ). The larger the  $n$ , the more information we have collected about the coefficients. However, even if the experiments are carried out with great care, they unavoidably will contain some error. The question remains: how may we estimate judiciously the coefficients  $c_1, \dots, c_k$  using the collected information about  $y = f_c(x)$ ? What is the best fit?

A very common method to respond to this question is known as the *method of least-squares*. The idea is simple: per realization of the experiment, we measure the fitting error by the distance from the real number  $f_c(x_1, x_2, \dots, x_m)$  and the observed value of  $y$ . The best fit for the distance though will also lead to a best fit for the square of the distance. To avoid absolute values, we

change our viewpoint slightly and measure the fitting error by  $(f_{c_1, \dots, c_k}(x_1, x_2, \dots, x_m) - y)^2$ . The error function that considers all the information obtained from the  $n$  experiments is then

$$E(c_1, \dots, c_k) = \sum_{i=1}^n ((f_{c_1, \dots, c_k}(x_{i1}, x_{i2}, \dots, x_{im}) - y_i)^2,$$

where  $(x_{i1}, x_{i2}, \dots, x_{im})$  and  $y_i$  are the  $i^{\text{th}}$  input for  $x$  and value of the measurement for  $y$ .

This turns out to be a function of  $c = (c_1, \dots, c_k)$ . Mathematically, our best fit problem is now reduced to finding the value of  $c$  which produces a minimum for this error function  $E$ . The details of how this can be done depends intrinsically upon the assumed form of the function  $f$ , and its relation to the parameters  $c_1, \dots, c_k$ .

We should remark that in some cases, we might want to use  $E(c_1, \dots, c_k)/n$ , the *mean-squared error*. This will not change the answer we get (because the minimum occurs at the same value of  $c$ ), but does allow us to compare sets of data with different numbers of data points.

### 3 Fitting a line to data

Suppose that we are in the situation described above, with  $m = 1$ , and that the function  $f$  is a polynomial of degree 1. Hence, the number of parameters is two, and for convenience we write

$$f(x) = mx + b,$$

in order to interpret the parameters  $m$  and  $b$  as the slope and  $y$ -intercept of the graph of  $f$ . We have data points  $(x_1, y_1), \dots, (x_n, y_n)$ . Our problem is to find the values of  $m$  and  $b$  will best fit this data. Or to put it simply and geometrically, what is the straight line that best fit the information contained in the data  $(x_1, y_1), \dots, (x_n, y_n)$ ?

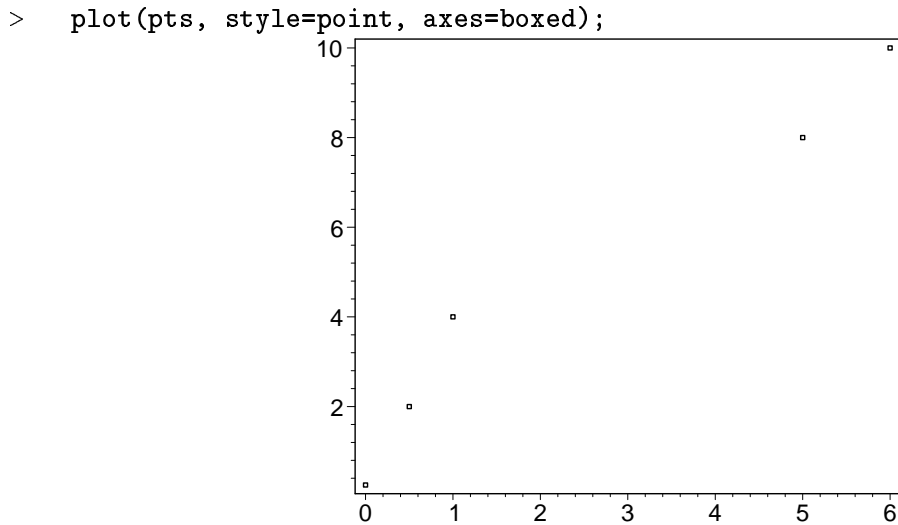
The contribution to the error function coming from the  $i^{\text{th}}$  piece of data,  $(x_i, y_i)$ , is given by  $(mx_i + b - y_i)^2$ , and therefore, the total error is

$$E(m, b) = \sum_{i=1}^n (mx_i + b - y_i)^2.$$

The problem then becomes that of searching for the values of  $m$  and  $b$  where the error  $E(m, b)$  achieves a minimum. Let us solve this problem using Maple.

First, let us type in some data to try to fit our line to:

```
> pts := [ [0,0.25], [1/2,2], [1,4], [5,8], [6,10] ]:
```



In order to solve the problem at hand, we could use Maple's built-in regression package to find the answer, but that would not help much in explaining what is going on. However, let us do this anyway.

We can use the `LeastSquares` command from the `CurveFitting` package<sup>3</sup>

```
> CurveFitting[LeastSquares](pts, x);
```

$$1.271370968 + 1.431451613x$$

This says that the best values of  $m$  and  $b$  are  $m = 1.431451613$  and  $b = 1.271370968$ , respectively. Now let's go through the steps ourselves, to better understand the process.

First, we define the error function according to the data points to fit that we have. Notice again that this function is the square of the distance from the line to the data:

```
> E := (m,b,pts)-> sum( ((m*pts[i][1]+b)-pts[i][2])^2, i=1..5);
```

$$E := (m, b, pts) \rightarrow \sum_{i=1}^5 (m \text{pts}_{i_1} + b - \text{pts}_{i_2})^2$$

For example, had we guessed that the line  $y = 2x + 1$  is the answer to our problem, we could compute the distance:

---

<sup>3</sup>As we noted before, the `CurveFitting` package was introduced in Maple 7. Earlier versions of Maple (as well as Maple 7) can accomplish the same thing using the `fit` command from the `stats` package, as follows.

```
with(stats):
fit[leastsquare][[x,y]]([[ pts[i][1] $i=1..5], [ pts[i][2] $i=1..5]]);
```

As you can see, the syntax is a little different.

```
> E(2,1,pts);
```

19.5625

However, decreasing  $m$  and increasing  $b$  produces a smaller value of  $E$ , so that guess cannot be the best fit:

```
> E(1.5,1.2,pts);
```

3.125000000

We want to minimize the error function. Since minima of differentiable functions occur at critical points, we begin by solving for the values of  $(m, b)$  which annihilate the two partial derivatives. Notice that Maple happily computes partial derivatives:

```
> diff( E(m,b,pts), m);
```

$$\frac{249}{2}m + 25b - 210$$

We may then ask Maple to solve the system of equations obtained by equating the two partial derivatives by

```
> solve( { diff( E(m,b,pts), m)=0, diff( E(m,b,pts), b)=0 }, {m,b});
```

$\{m = 1.431451613, b = 1.271370968\}$

Maple thus finds only one solution, and not surprisingly, it coincides with the solution found using the built-in version.

If we now want to use this values of  $m$  and  $b$  without retyping it, one way to do that is to use the command `assign(%)` to let  $b$  and  $m$  be given these constant values.<sup>4</sup>

```
> assign(%)
```

Maple executes this assignment silently, without reporting the assignment. However, now both  $m$  and  $b$  have the assigned values:

```
> m;
```

1.431451613

---

<sup>4</sup>We are using Maple's *ditto operator* % to refer to the result of the previous command. In versions of Maple prior to release 5 of Maple V, the double quote " was used for this operator instead.



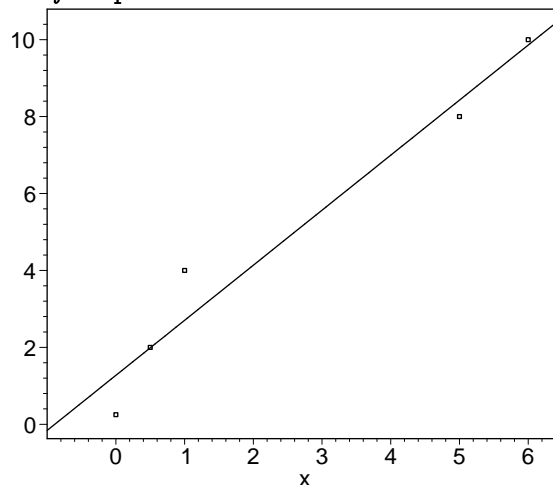
```
> E(m,b,pts);
```

2.929334677

The value of  $E(m, b, pts)$  calculated above, is the value of  $E$  when  $(m, b)$  is the critical point found earlier.

We may visualize how good our fit is, using `plot` to display the data points and the fit line on the same graph. We load the `plots` package, so that we can use `display`:

```
> with(plots):
  display({plot(m*x+b, x=-1..6.5, axes=boxed),
          plot(pts, style=point)}):
```



Notice that at this point,  $m$  and  $b$  have constant values previously assigned to them, the values producing the minimum of  $E$  up to 9 decimal places. Let us try to verify that our solution is indeed correct, via an argument which is independent of the one described above.

Geometrically, the error function  $E$  is the square of the Euclidean distance between  $m\vec{x} + \vec{b}$  and the vector  $\vec{y}$ , where  $\vec{x} = (x_1, x_2, \dots, x_n)$ ,  $\vec{b} = (b, b, \dots, b)$  and  $\vec{y} = (y_1, \dots, y_n)$ . Clearly then, the best fit is achieved when the coefficients  $m$  and  $b$  are chosen so that  $m\vec{x} + \vec{b}$  is the orthogonal projection of the vector  $\vec{y}$  onto the subspace of  $\mathbb{R}^n$  spanned by  $\vec{x}$  and  $(1, 1, \dots, 1)$ . That is to say, the best fit occurs when the vectors  $\vec{y} - (m\vec{x} + \vec{b})$  and  $m\vec{x} + \vec{b}$  are perpendicular to each other. Let us verify this:

```
> sum( (pts[i][2]-m*pts[i][1]-b)*(m*pts[i][1]+b), i=1..5);
```

$-0.16 \cdot 10^{-7}$

The resulting inner product is not quite zero, but the first non-zero digit of its decimal expansion occurs in the eighth decimal place. This is because the calculations were only done with finite precision, and the discrepancy above is attributable to round-off errors..

As a final remark, we observe that if you now execute the command `diff( E(m,b,pts),m )`, you will obtain an error. This is due to the fact that  $m$  and  $b$  are constants at this point, and not independent variables. If you would like to continue making use of the function  $E(m, b, pts)$  for general  $m$  and  $b$ , you will have to unassign the values of  $m$  and  $b$  first:

```
> m := 'm':
   b := 'b':
```

## 4 Fitting a cubic to data

Let us now try to fit a cubic polynomial to some data. We do so directly, without using Maple's built-in fitting package.

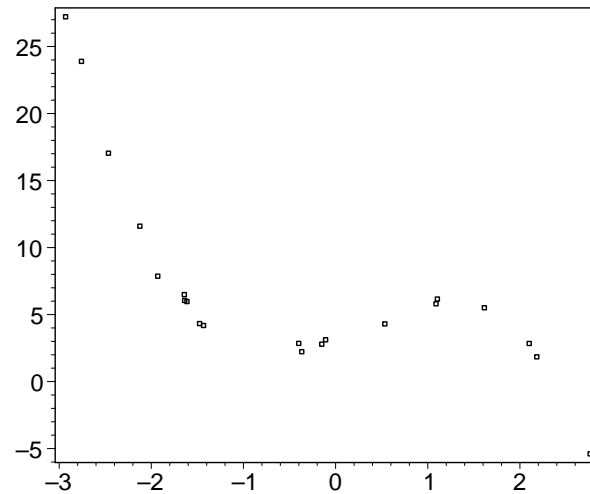
We randomly generate a set of 21 data points to be fitted to a cubic using a Maple program. (Don't worry about how this program works just now. We'll get into those details later).

```
> Seed:=readlib(randomize()):
   with(stats):
   myrand := (a,b) -> evalf(a+(b-a)*(rand()/10^12)):
   points_wanted := 21:

> cubic_pts := proc()
   local fuzz, a, b, c, d,s,x,i;
   a:= myrand(-3,3);
   b:= myrand(-3,3);
   c:= myrand(-3,3);
   d:= myrand(-10,10);
   if (myrand(-10,10) > 0) then s:=1 else s:=-1 fi;
   fuzz := [ random[normald[0,0.5]](points_wanted)];
   x     := [ random[uniform[-3,3]](points_wanted)];
   RETURN([seq([x[i],s*(a-x[i])*(b-x[i])*(c-x[i])+d+fuzz[i]],
               i=1..points_wanted)]);
end:
```

Now we can put the data to be fitted into a list, and visualize the result with `plot`:

```
> data:=cubic_pts():
   plot(data, style=point, axes=boxed);
```



Now we have a set of 21 data points which are to be fitted to a cubic polynomial. A polynomial function of degree  $n$  is determined by  $n + 1$  coefficients. So we define

```
> cub := (x,a,b,c,d)->a*x^3+b*x^2+c*x+d;
```

$$\text{cub} := (x, a, b, c, d) \rightarrow ax^3 + bx^2 + cx + d$$

and then define the error function:

```
> E:=(a,b,c,d,data)->sum((cub(data[i][1],a,b,c,d)-data[i][2])^2,i=1..nops(data));
```

$$E := (a, b, c, d, data) \rightarrow \sum_{i=1}^{\text{nops}(data)} (\text{cub}(data_{i1}, a, b, c, d) - data_{i2})^2$$

We have four parameters to determine, the coefficients of the function `cub`. We find its values by

```
> assign(solve({diff(E(a,b,c,d,data),a)=0,
                diff(E(a,b,c,d,data),b)=0,
                diff(E(a,b,c,d,data),c)=0,
                diff(E(a,b,c,d,data),d)=0},
                {a,b,c,d}));
```

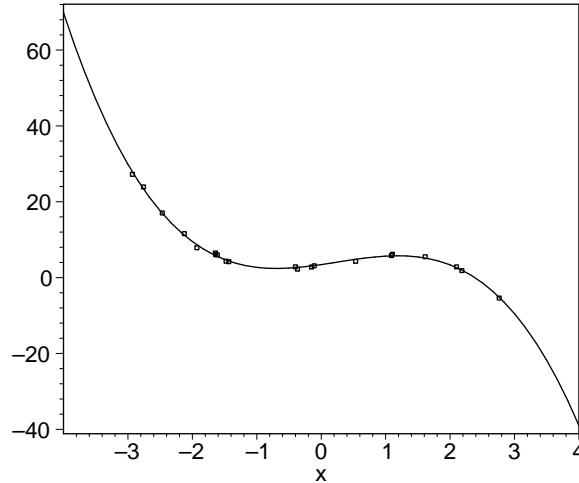
Of course, we do not see the answer, but the resulting values have already being assigned to the coefficients  $a$ ,  $b$ ,  $c$ ,  $d$ :

```
> cub(x,a,b,c,d);
```

$$-1.006141915x^3 + .7505869176x^2 + 2.487673553x + 3.430571969$$

Finally, let's make a picture of the result, showing both the data and the fitted curve:

```
> with(plots):
cplot := plot(cub(x,a,b,c,d), x=-4..4, axes=boxed):
pplot := plot(data, style=point):
display({cplot,pplot});
```



For comparison, let us find the answer using Maple's built-in statistical package:

```
> with(CurveFitting):
LeastSquares(data,x,curve=A*x^3+B*x^2+C*x+D);
```

$$-1.006141912x^3 + .7505869193x^2 + 2.487673537x + 3.430571965$$

For this particular version of least square fitting, we may once again interpret the error function  $E$  geometrically as the square of the Euclidean distance between  $a\vec{x}_3 + b\vec{x}_2 + c\vec{x}_1 + \vec{d}$  and the vector  $\vec{y}$ , where  $\vec{x}_j = (x_1^j, x_2^j, \dots, x_n^j)$ ,  $\vec{d} = d(1, 1, \dots, 1)$  and  $\vec{y} = (y_1, \dots, y_n)$ . The best fit is thus achieved when the coefficients  $a$ ,  $b$ ,  $c$  and  $d$  are chosen so that  $a\vec{x}_3 + b\vec{x}_2 + c\vec{x}_1 + \vec{d}$  is the orthogonal projection of the vector  $\vec{y}$  onto the subspace of  $\mathbb{R}^n$  spanned by  $\vec{x}_3$ ,  $\vec{x}_2$ ,  $\vec{x}_1$  and  $(1, 1, \dots, 1)$ . That is to say, the best fit occurs when the vectors  $\vec{y} - (a\vec{x}_3 + b\vec{x}_2 + c\vec{x}_1 + \vec{d})$  and  $a\vec{x}_3 + b\vec{x}_2 + c\vec{x}_1 + \vec{d}$  are perpendicular. You should verify that this is the case for the solution given above.

In both cases, you should consider why there is only one critical point for the error function, and why it is of necessity a global minimum.

## 5 Fitting other types of funtions

In the previous constructions, we dealt with functions that depended linearly upon the parameters to be fitted. We found them by solving a linear system of equations, a relatively easy task.

However, as explained at the beginning of this chapter, the same scheme works equally

well for any function. In the general case, though, the resulting system that we must solve to find the best fit could depend non-linearly on the parameters, and the mere existence of solutions to such systems is not a trivial problem to settle. Take for instance, a function such as  $y = \sum_{i=1}^n \sin(m_i x)$ . Assuming measurements  $(x_1, y_1), \dots, (x_k, y_k)$ , we may take as our error the function

$$E(m_1, \dots, m_n) = \sum_{j=1}^k (y_j - \sum_{i=1}^n \sin(m_i x_j))^2.$$

Its partial derivatives are quite easy to calculate, but it is far from clear if there are values of  $m_1, \dots, m_n$  where they all vanish simultaneously. For a particular set of values, we may ask Maple to solve the resulting system and get absolutely nothing. This either indicates an inability of Maple to handle such a system, or worse yet, the fact that such a system has no solution at all. Even when the latter happens, Maple itself might not be able to tell us so.

Sometimes, even if the function does not depend linearly on the coefficients, it can be transformed to one which does. For example, if our data points  $\{(x_i, y_i)\}$  were believed to approximate an exponential function of the form  $y = ae^{kx}$ , then setting

$$E(a, k) = \sum_{i=1}^m (y_i - ae^{kx_i})^2.$$

would require us to solve the system

$$\sum_{i=1}^m (y_i - ae^{kx_i})e^{kx_i} = 0 \quad \sum_{i=1}^m (y_i - ae^{kx_i})ax_i e^{kx_i} = 0.$$

While this isn't impossible, it is much more straightforward to make a change of variables.

Assuming the  $y_i$  are all positive (which is reasonable since we believe the data is exponential), we can let  $z_i = \ln y_i$ . Then we are trying to fit a function  $z = \ln(ae^{kx})$ , which reduces to the line  $z = \ln a + kx$ . This is familiar territory, and we can just proceed as before.

## 6 Fitting a circle

In a slightly different vein, suppose that the data points  $\{(x_i, y_i)\}$  lie near a circle of unknown center and radius. Note that if we did know the center, this problem would be very simple: the desired radius would be the average distance from the points to the center. If you don't think very hard, you might suspect that the desired center would be very near the center of mass of the points  $\{(x_i, y_i)\}$ , but this will only be the case if the points are evenly distributed around the circle.

Instead, we can come up with a function that measures the error between an arbitrary circle and our data points.

Recall that the general equation of a circle centered at  $(a, b)$  with a radius of  $r$  is

$$(x - a)^2 + (y - b)^2 = r^2.$$

Unlike the previous cases, there is no independent variable. (Note that we *could* try to fit  $y_i \simeq b \pm \sqrt{r^2 - (x_i - a)^2}$ , but not only would the resulting equations be messy, this would bias things very badly; do you see why?). Nevertheless, we press on.

One reasonable measure of the distance between the points and a circle is the “area difference”, that is

$$E(a, b, r) = \sum_i ((x_i - a)^2 + (y_i - b)^2 - r^2)^2.$$

One difficulty with this is that  $E$  is not quadratic in  $a$ ,  $b$ , and  $r$ . However, Maple is able to solve the resulting equations anyway.

We assume here that there is a routine called `circle_pts` which gives us points which lie near a circle of unknown radius and center. The one we use here is similar to that in §4, and is defined in `lsq_data.txt`.

```
> cpts:=circle_pts():
   epsilon:=(pt,a,b,r) -> (pt[1]-a)^2 + (pt[2]-b)^2 -r^2;
```

$$\varepsilon := (pt, a, b, r) \rightarrow (pt_1 - a)^2 + (pt_2 - b)^2 - r^2$$

```
> E:= (a,b,r,pts) -> sum( epsilon(pts[i],a,b,r)^2, i=1..nops(pts));
```

$$E := (a, b, r, pts) \rightarrow \sum_{i=1}^{\text{nops}(pts)} \varepsilon(pts_i, a, b, r)^2$$

```
> sol:= {solve({diff(E(a,b,r,cpts),a)=0,
               diff(E(a,b,r,cpts),b)=0,
               diff(E(a,b,r,cpts),r)=0},
             {a,b,r})};
```

```
sol := [ {r = 0., b = 8.598845417 - 5.760818468 I, a = 10.74842676 - 3.044017219 I},
         {r = 0., a = 10.74842676 + 3.044017219 I, b = 8.598845417 + 5.760818468 I},
         {r = 0., b = 9.123213697, a = 8.304601815},
         {r = 0., a = 10.50158731 + 4.256560324 I, b = 11.57538622 - 5.447431926 I},
         {b = 11.57538622 + 5.447431926 I, r = 0., a = 10.50158731 - 4.256560324 I},
         {a = 8.224446708, r = -4.468796472, b = 8.538288341},
         {a = 8.224446708, b = 8.538288341, r = 4.468796472}]
```

Here we find a number of critical points for the function  $E$ , but from physical considerations, the only reasonable choice is the circle for which  $r > 0$ .

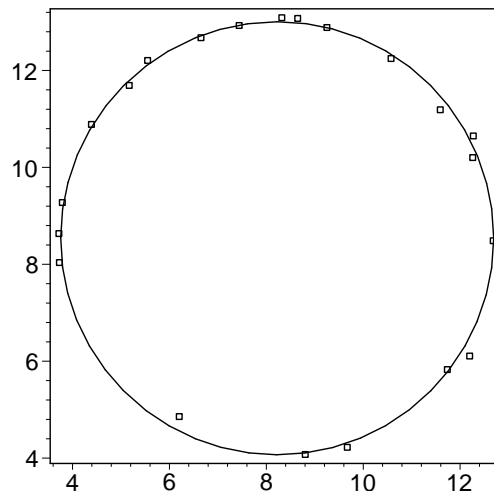
This is the seventh entry in the above list, so we could refer to it just by `sol[7]`. However, with a different set of data, it might be the second solution, or the fourth, etc. To ensure that we always get the right one, no matter what the order is, we can use Maple's `select` command to pull out the one with  $r > 0$ . This looks more daunting than it really is: we just define a function that returns `true` if that  $r > 0$  for that solution, and `false` if not. Then `select` returns only those elements for which the function is true. We need to use `op`, because we have a list of sets, and we just want the one answer.

```
> goodsol:=op(select(s->if (subs(s,r)>0) then true else false fi,sol));
```

```
goodsol := {a = 8.224446708, b = 8.538288341, r = 4.468796472}
```

In order to see the fit, we plot the points and the circle. Note that we represent the circle parametrically, and use `subs` to substitute the desired solution.

```
> display(
  plot(cpts,style=point,scaling=constrained,axes=boxed),
  plot(subs(goodsol,[a+r*cos(t),b+r*sin(t),t=0..2*Pi]))
);
```



The fact that there is only one interesting critical point is no accident, however. If we let  $k = a^2 + b^2 - r^2$ , then the resulting error function  $H(a, b, k) = E(a, b, \sqrt{a^2 + b^2 - k})$  is quadratic in  $a$ ,  $b$ , and  $k$ , and so we really only need to solve a linear system. It is easy to check that this new functional  $H$  still satisfies all the criteria we wanted (that is  $H(a, b, k) = 0$  if and only if all the data points lie on the circle  $C(a, b, k)$ ,  $H$  is non-negative, and it is smooth), so fitting a circle becomes a linear problem after all. The reader should verify that, in fact, the unique minimum of  $H$  corresponds exactly to the critical points the original function  $E$  for which  $r \neq 0$ .

The interested reader might want to consider a similar approach to fitting other conic

sections, such as an ellipse or a hyperbola, to given data. Do you expect to be able to make the problem linear, as in the case of the circle?

## 7 Robust fitting

Let us return momentarily to the problem discussed in §3: given some data  $(x_1, y_1), \dots, (x_n, y_n)$ , we found the best line that fits it. This was accomplished by measuring the square of the *vertical* distance from each data point and the line  $y = mx + b$ , which led us to consider the error function

$$E(m, b) = \sum_{i=1}^n (mx_i + b - y_i)^2.$$

The problem was solved by finding the values of  $m$  and  $b$  where  $E$  achieves its minimum.

However, the square of the vertical distance is only one of many reasonable ways of measuring the distance from the data points to the line  $mx + b$ . For example, it is perfectly reasonable to consider instead the expression  $1 + (mx_i + b - y_i)^2$ . If  $(x_i, y_i)$  is on the line, this value would be equal to 1, and its logarithm would then be zero. Thus, the function

$$E(m, b) = \sum_{i=1}^n \ln(1 + (mx_i + b - y_i)^2),$$

is also a reasonable way of measuring the distance from the data points to the line  $y = mx + b$ , and its minimum would produce a best fit line in this other sense. Similarly, We may argue that the function

$$E(m, b) = \sum_{i=1}^n |mx_i + b - y_i|.$$

is another reasonable alternative, though this time  $E$  is not even differentiable in the region of interest (if a data point  $(x_i, y_i)$  happens to be exactly on the line,  $mx_i + b - y_i = 0$  and the absolute value is not differentiable at 0).

There is a priori no particular reason to prefer one error function over another. And for each choice of such we could end up with a problem whose solution exists and produces a line that best fit the data. Then we could compute the value of the error function for the best fit line, value that depends upon the choice of function made. A *canonical* way of choosing the best error function could be that for which this number is the smallest. But finding such a function is quite a difficult, if not impossible, problem to solve. If we limit our attention to *linear estimators*, in a sense to be explained in the last section of this chapter, the solution obtained in §3 is optimal. In here, we pursue further the problem of finding the best fit according to the two error functions introduced above.

Notice that if  $\epsilon_i = mx_i + b - y_i$ , data with large errors  $\epsilon_i$  exert a huge influence on the original error function of §3 (because we use  $\sum \epsilon_i^2$ ). The two error functions mentioned above



grow linearly (or sub-linearly) with  $|\epsilon_i|$ , causing the tails to count much less heavily. We will see this when comparing the different results.

In order to solve the problem now at hand, we first generate some “noisy” data. We load the routines in `lsq_data.txt`:

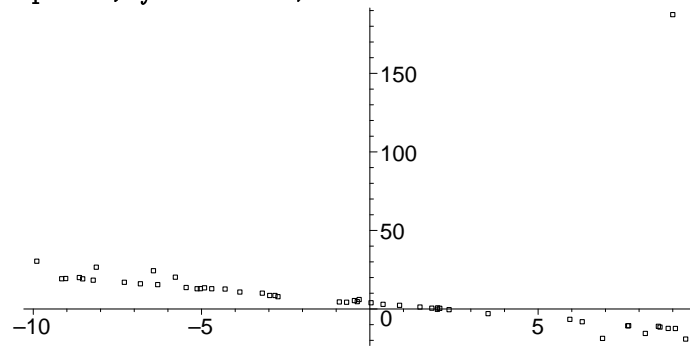
```
> read('lsq_data.txt');
```

```
defined line_pts(), bad_line_pts(), quadratic_pts(), cubic_pts(), and circle_pts()
```

The data and its graphical visualization can then be obtained by:

```
> pts:=bad_line_pts():
```

```
> plot(pts,style=point,symbol=box);
```



You should notice that while most of the data lies fairly near a line, there is one point which is extremely far away from the rest of the data.

In order to compare results, let us first find the line given by least squares. First, we define a function  $\epsilon$ , which gives the signed vertical distance between a point `pt` and a line of slope `m` and intercept `b`

```
> epsilon := (pt,m,b) -> (m*pt[1]+b-pt[2]);
```

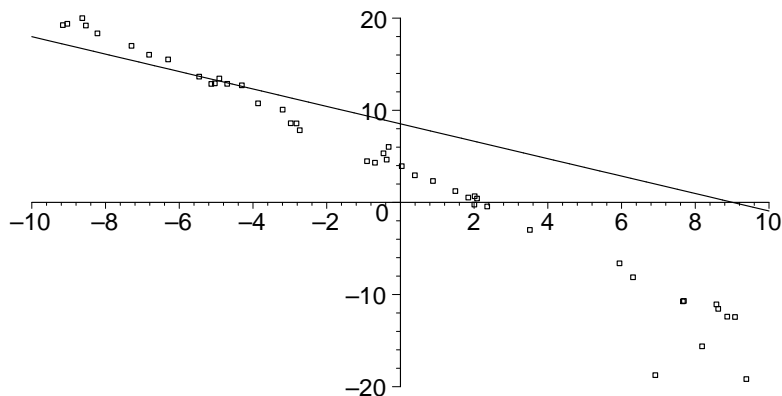
$$\epsilon := (pt, m, b) \rightarrow m pt_1 + b - pt_2$$

Now we compute the regular least squares fit.

```
> H:=(m,b,pts)->sum(epsilon(pts[i],m,b))^2,i=1..nops(pts):
sol:=solve({diff(H(m,b,pts),m)=0, diff(H(m,b,pts),b)=0},{m,b});
```

$$sol := \{m = -.9455490512, b = 8.528125863\}$$

```
> display({plot(pts,style=point,symbol=box),plot(subs(sol,m*x+b),x=-10..10)},
view=-20..20);
```



In the above, we restricted our attention to the region  $-20 < y < 20$ , where the line and the majority of the data lies. As you can see, the fit to the majority of the data is very poor, because of the influence exerted by the anomalous data point.

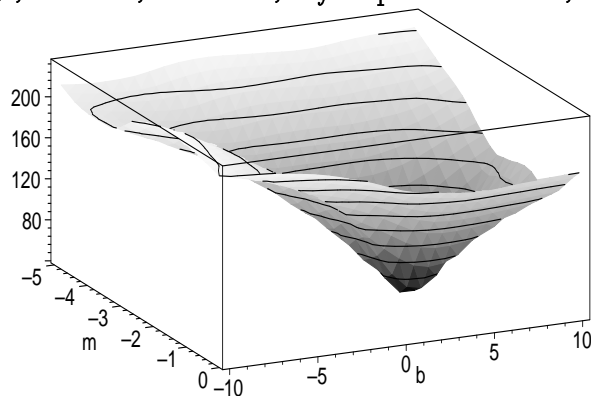
Now, let's try again, this time minimizing  $\ln(1 + \epsilon_i^2/2)$ , which behaves like  $\epsilon^2$  for small errors, but grows very slowly for large ones.

```
> R:=(m,b,pts)->sum(ln(1+epsilon(pts[i],m,b)^2/2),i=1..nops(pts));
```

$$R := (m, b, pts) \rightarrow \sum_{i=1}^{\text{nops}(pts)} \ln\left(1 + \frac{1}{2}\epsilon(pts_i, m, b)^2\right)$$

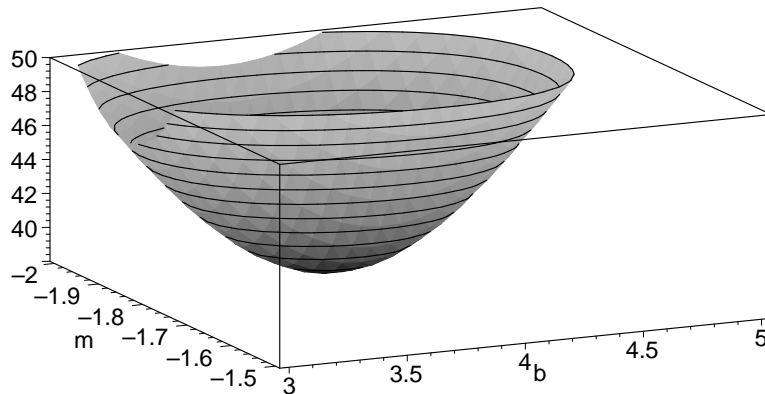
Here's what the functional we want to minimize looks like:

```
> plot3d(R(m,b,pts),m=-5..0,b=-20..0,style=patchcontour, axes=boxed);
```



The equations we want to solve are not linear. In fact, they're messy enough that Maple needs some help to find the minimum. By examining the plot and zooming in on the minimum, we can choose an appropriate region.

```
> plot3d(R(m,b,pts),m=-2..-1.5,b=3..5,view=38..50, axes=boxed);
```



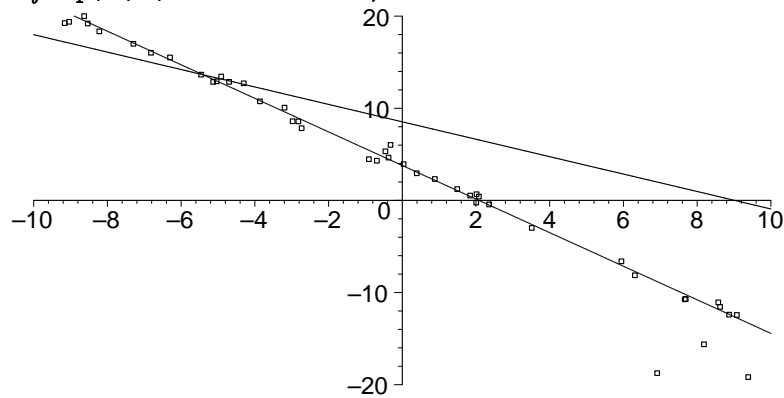
Once we have an appropriate region, we can ask Maple to look for the solution numerically using `fsolve`.

```
> rs:=fsolve( {diff(R(m,b,pts),m)=0, diff(R(m,b,pts),b)=0},
               {m,b},m=-2..-1.5, b=3..5);
```

```
rs := {m = -1.822071903, b = 3.789519722}
```

And here we can compare the two lines found.

```
> p := plot(pts,style=point):
  l := plot(subs(sol,m*x+b),x=-10..10):
  s := plot(subs(rs,m*x+b),x=-10..10,color=blue):
  plots[display](p,l,s,view=-20..20);
```



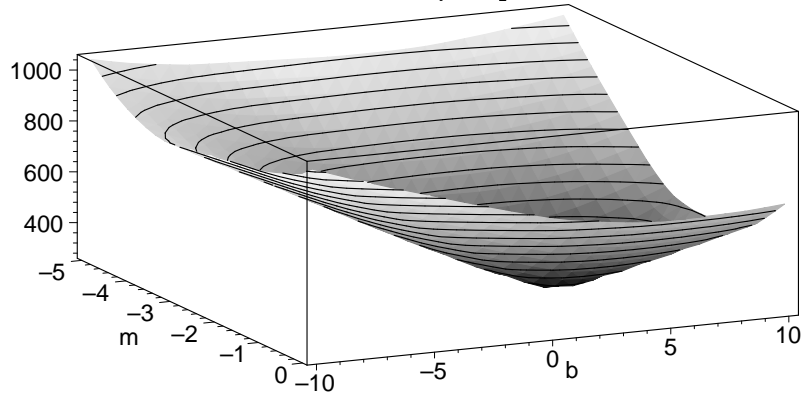
Let us now try to minimize  $\sum |\epsilon_i|$ . Since this function is not even differentiable near  $\epsilon_i = 0$ , we have to work harder to get less.

```
> A:=(m,b,pts)->sum(abs(epsilon(pts[i],m,b)),i=1..nops(pts));
```

$$A := (m, b, pts) \rightarrow \sum_{i=1}^{\text{nops}(pts)} |\epsilon(pts_i, m, b)|$$

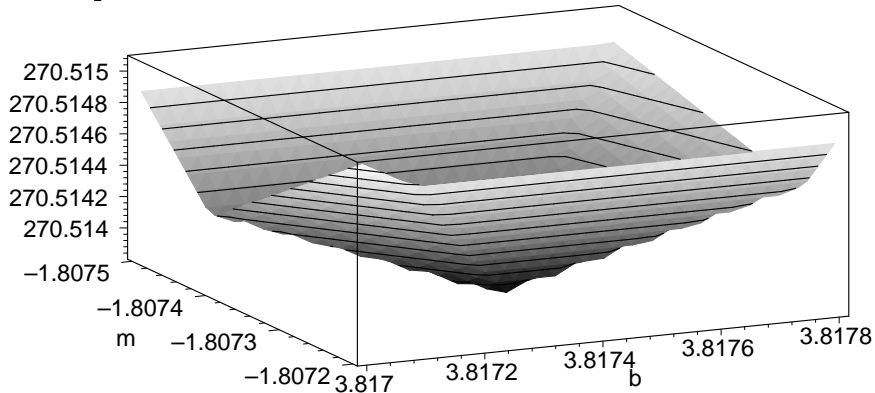
Displaying the graph of this function gives us an idea about where its minimum lies:

```
> plot3d(A(m,b,pts),m=-5..0,b=-10..10,style=patchcontour, axes=boxed);
```



But coercing Maple into finding a numerical approximation to it is not so easy. There *are* reliable and efficient ways to determine the location of the minimum to any precision (at least given some hints), but they are inappropriate for the scope of this chapter, and will not be considered here. Instead we may zero in on the minimum by looking at where the lowest point is and repeatedly restricting the domain. After several iterations, we get the following.

```
> plot3d(A(m,b,pts),m=-1.0875..-1.0872,b=3.817..3.8178, axes=boxed);
```



This leads us to a choice of  $m = -1.80733$ ,  $b = 3.8174$  as our “best” line. This is almost (but not quite) the same line found using our distance  $R(m, b)$ .

## 8 A nod toward statistics

In the problem we have treated so far, there have been a distinction made between the  $x$  and  $y$  coordinate of a data point  $(x_i, y_i)$ : the  $x$  coordinate is thought as the predictor variable, while the  $y$  coordinate is the predicted value. The distinction is significant. If, for example, we think of  $x$  as a function of  $y$ , the best line that fits the data using the method of least square is, in general, different than the corresponding line when thinking of  $y$  as a function of  $x$ . Indeed, if we use the same data as in §3, the best line  $x = ny + c$  that fits it is given by

$x = 0.6677953348y - 0.738807374$ . Solving in terms of  $y$  yields  $y = 1.497464789x + 1.106338028$ ; this differs significantly from the line  $y = 1.431451613x + 1.271370968$  found in §3.

In this case, we use as error the square of the *horizontal* distance, and it is somewhat perturbing that this similarly reasonable approach leads to a best fit that differs from the one obtained when employing vertical distances.

The situation can be made even worse if neither  $x$  nor  $y$  is a predictor variable. In this case, we want to minimize the shortest distance to a line  $y = mx + b$  rather than the vertical (or horizontal) distance. The resulting equations will not be linear, nor can they be made linear. However, Maple will be able to find the critical points with no trouble. There will always be at least two (there may be a third with a huge slope and intercept)—one is the minimum, and the other is a saddle point. It is worth thinking for a few minutes about the geometric interpretation of the saddle point in terms of the problem at hand.

In practice, the perturbing fact that different but reasonable error functions lead to best fits that are different is resolved by knowing what part of the data is the input and what part is predicted from it. This happens often enough, though not always. We thus can make a good choice for  $E$  and solve a minimization problem. That settled, we are still left with the problem of demonstrating why our choice for  $E$  is a good one.

It is rather easy to write down the solution to the problem in §3: if

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i,$$

then

$$\hat{m} = \frac{n \sum_{i=1}^n x_i y_i - \left( \sum_{j=1}^n x_j \right) \left( \sum_{j=1}^n y_j \right)}{n \sum_{i=1}^n x_i x_i + \left( \sum_{j=1}^n x_j \right)^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \hat{b} = \bar{y} - m\bar{x}.$$

If  $y_i$  are assumed to be the values of random variables  $Y_i$  which depend *linearly* upon the  $x_i$ ,

$$Y_i = mx_i + b + \varepsilon_i,$$

with errors  $\varepsilon_i$  that are independent from each other, are zero on average and have some fixed variance, then the value of  $m$  given above is the *best* estimator of the slope among all those linear estimators that are unbiased. “Best” here is measured by calculating the deviation from the mean, and this best estimator is the one that produces the smallest such deviation.

This conclusion follows by merely making assumptions about the inner products of the data points  $(x_1, \dots, x_n)$  and  $(y_1, \dots, y_n)$ . Statisticians often would like to answer questions such as the degree of accuracy of the estimated value of  $m$  and  $b$ . For that one would have to assume more about the probability distribution of the error variables  $Y_i$ . A typical situation is to assume that the  $\varepsilon_i$  above are normally distributed, with mean 0 and variance  $\sigma^2$ . Under these assumptions, the values of  $m$  and  $b$  given above are the so-called *maximum likelihood estimators* for these

two parameters, and there is yet one such estimator for the variance  $\sigma^2$ . But, since we assumed more, we can also say more. The estimators  $\hat{m}$  and  $\hat{b}$  are normally distributed and, for example, the mean of  $\hat{m}$  is  $m$  and its variance is  $\sigma^2 / \sum_{i=1}^n (x_i - \bar{x})^2$ . With this knowledge, one may embark into determining the confidence we could have on our estimated value for the parameters. We do not do so in here, but want to plant the idea in the interested reader, whom we refer to books on the subject.

# The Art of Phugoid

In this chapter, we will explore some aspects of mathematical model of glider flight. This model is called Lanchester's *phugoid theory*<sup>3.1</sup>, developed by Frederick Lanchester [[Lan](#)] at beginning of the twentieth century. While this model has its drawbacks, it still is used today to explain oscillations and stalls in airplane flight.

Since this model, like many models, is a set of differential equations, we will need some results which are typically covered in a course on differential equations. We will cover the relevant material briefly here, but readers needing a more in-depth treatment are encouraged to look in a text on ordinary differential equations, such as [[BDH](#)], [[HW](#)], or [[EP](#)]. Of these, the approach in [[BDH](#)] is perhaps closest to the one presented here.

---

## Footnotes

... theory<sup>3.1</sup>

Folklore has it that this name was chosen by Lanchester because he wanted a classically-based name for his new theory of oscillations occurring during flight. Since the Greek root *phug* ( $\phi\upsilon\gamma\eta$ ,

pronounced ``fyoog") as well as the Latin root *fug*- both correspond to the English word ``flight", he decided on the name ``phugoid". Unfortunately, in both Greek and Latin, this means ``flight" as in ``run away" instead of what birds and airplanes do-- the same root gives rise to the words ``fugitive" and ``centrifuge". The Latin for flight in appropriate sense is *volatus*; the Greek word is *potê* ( $\pi\omicron\tau\eta$ ).

---

## Subsections

- [The Phugoid model](#)
  - [What do solutions look like?](#)
  - [Existence of Solutions](#)
  - [Numerical Methods](#)
    - [Euler's method](#)
    - [Numerical Solutions, Numerical Integration, and Runge-Kutta](#)
  - [Seeing the flight path](#)
  - [Fixed Point Analysis](#)
    - [Linear Systems of ODEs](#)
    - [Fixed Points for the Glider](#)
  - [Qualitative Classification of Solutions](#)
  - [Dealing with the Singularity](#)
-

[next](#) [up](#) [previous](#)

**Next:** [The Phugoid model](#) **Up:** [Using Maple](#) **Previous:** [A nod toward statistics](#)

*Translated from LaTeX by Scott Sutherland  
2002-08-29*



# Chapter 3

## The Art of Phugoid

In this chapter, we will explore some aspects of mathematical model of glider flight. This model is called Lanchester’s *phugoid theory*<sup>1</sup>, developed by Frederick Lanchester [Lan] at beginning of the twentieth century. While this model has its drawbacks, it still is used today to explain oscillations and stalls in airplane flight.

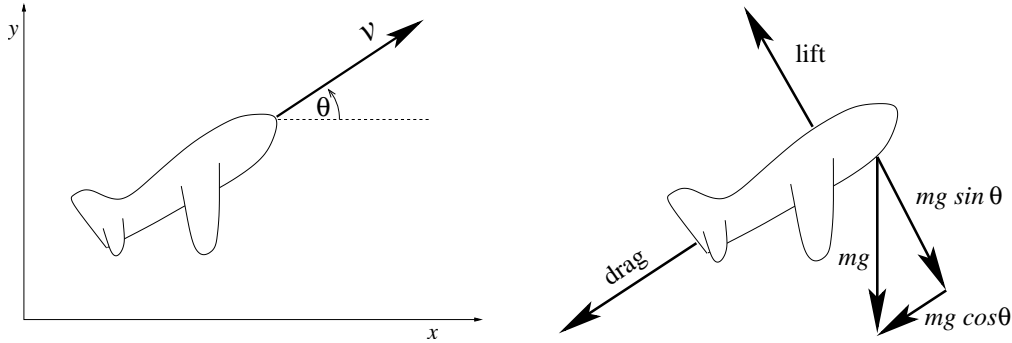
Since this model, like many models, is a set of differential equations, we will need some results which are typically covered in a course on differential equations. We will cover the relevant material briefly here, but readers needing a more in-depth treatment are encouraged to look in a text on ordinary differential equations, such as [BDH], [HW], or [EP]. Of these, the approach in [BDH] is perhaps closest to the one presented here.

### 1 The Phugoid model

The Phugoid model is a system of two nonlinear differential equations in a frame of reference relative to the plane. Let  $v(t)$  be the speed the plane is moving forward at time  $t$ , and  $\theta(t)$  be the angle the nose makes with the horizontal. As is common, we will suppress the functional notation and just write  $v$  when we mean  $v(t)$ , but it is important to remember that  $v$  and  $\theta$  are functions of time.

---

<sup>1</sup>Folklore has it that this name was chosen by Lanchester because he wanted a classically-based name for his new theory of oscillations occurring during flight. Since the Greek root *phug* ( $\phi\nu\gamma\eta$ , pronounced “fyoog”) as well as the Latin root *fug-* both correspond to the English word “flight”, he decided on the name “phugoid”. Unfortunately, in both Greek and Latin, this means “flight” as in “run away” instead of what birds and airplanes do—the same root gives rise to the words “fugitive” and “centrifuge”. The Latin for flight in appropriate sense is *volatus*; the Greek word is *potê* ( $\rho\omicron\tau\eta$ ).



If we apply Newton's second law of motion (force = mass  $\times$  acceleration) and examine the major forces acting on the plane, we see easily the force acting in the forward direction of the plane is

$$m \frac{dv}{dt} = -mg \sin \theta - \text{drag}.$$

This matches with our intuition: When  $\theta$  is negative, the nose is pointing down and the plane will accelerate due to gravity. When  $\theta > 0$ , the plane must fight against gravity.

In the normal direction, we have centripetal force, which is often expressed as  $mv^2/r$ , where  $r$  is the instantaneous radius of curvature. After noticing that that  $\frac{d\theta}{dt} = v/r$ , this can be expressed as  $v \frac{d\theta}{dt}$ , giving

$$mv \frac{d\theta}{dt} = -mg \cos \theta + \text{lift}.$$

Experiments show that both drag and lift are proportional to  $v^2$ , and we can choose our units to absorb most of the constants. Thus, the equations simplify to the system

$$\frac{dv}{dt} = -\sin \theta - Rv^2 \qquad \frac{d\theta}{dt} = \frac{v^2 - \cos \theta}{v}$$

which is what we will use henceforth. Note that we must always have  $v > 0$ .

It is also common to use the notation  $\dot{v}$  for  $\frac{dv}{dt}$  and  $\dot{\theta}$  for  $\frac{d\theta}{dt}$ . We will use these notations interchangeably.

## 2 What do solutions look like?

Now that we have some equations in hand, our first goal is to understand what the solutions look like. Note that the equations are nonlinear, and we cannot hope to find explicit formulae for the solutions—they don't exist except in very special cases.

Let's deal with the easiest case first: pretend that drag doesn't exist and set  $R = 0$ . Does this help?

The equations become

$$\frac{dv}{dt} = -\sin \theta \qquad \frac{d\theta}{dt} = \frac{v^2 - \cos \theta}{v}$$

Now, let's see what we can say about the solutions to this equation. First, notice that there must be many solutions: Any choice of initial conditions  $\theta(0) = \theta_0, v(0) = v_0$  gives rise to a different solution. (Informally, if the glider starts with a different initial angle and/or velocity, it will fly along a different path. If we say it this way, it seems so obvious we shouldn't need to say it at all.) Can we find *any* solutions at all?

There is one very easy solution to find. Notice that if we can find a choice of  $\theta$  and  $v$  so that  $\frac{dv}{dt} = 0 = \frac{d\theta}{dt}$ , then since their derivatives are zero,  $v(t)$  and  $\theta(t)$  must be constant functions. But solving this is simple:

$$0 = -\sin \theta \qquad 0 = \frac{v^2 - \cos \theta}{v}$$

is easily seen to be true when  $\theta = 0$  and  $v = 1$ . In fact,  $\theta$  can be any multiple of  $2\pi$ . So we have found the special solution

$$v(t) \equiv 1 \qquad \theta(t) \equiv 0.$$

This is a glider which always flies level, with a constant velocity of 1.

What about other solutions? Now we must work harder. There is no hope to find a general analytic expression for  $v(t)$  and  $\theta(t)$ , but that does not mean all is lost. We can get an idea of how other solutions behave by thinking about what it means to have a solution  $\varphi(t) = [\theta(t), v(t)]$ .

First, notice that such a solution  $\varphi$  can be thought of as a parametric curve in the  $(\theta, v)$ -plane which passes through the point  $\varphi(0) = [\theta_0, v_0]$ . At any point along this parametric curve, we can calculate its derivative:

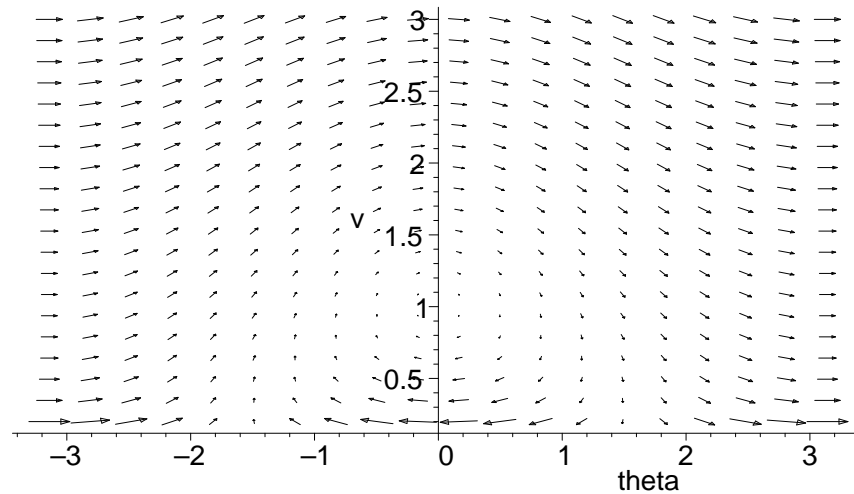
$$\varphi'(t) = \left[ \frac{d}{dt}\theta(t), \frac{d}{dt}v(t) \right] = [\dot{\theta}(t), \dot{v}(t)].$$

But we started out with expressions for  $\dot{\theta}$  and  $\dot{v}$ : our original differential equation. Because our differential equations only involve  $t$  implicitly via  $\theta$  and  $v$  (the system is *autonomous*), we need only know the values of  $\theta$  and  $v$  to know the tangent vector to any solution curve passing through that point.

We can get an idea of how the solutions behave by making a picture of the vector field corresponding to this: At each value for the pair  $(\theta, v)$ , we draw the vector  $\langle \dot{\theta}, \dot{v} \rangle$ , which only depends on  $(\theta, v)$ .

To get maple to do this, we can use the `fieldplot` command. Because the arrows in the `fieldplot` show not only direction but also relative magnitude, we shouldn't get too close to the  $v = 0$  axis or they will overwhelm the smaller arrows for which occur for larger  $v$ .

```
> plots[fieldplot]([ (v^2 - cos(theta))/v, -sin(theta) ],
                   theta=-Pi..Pi, v=0.2..3, arrows=slim);
```



We can certainly see that near the constant solution  $\{\theta = 0, v = 1\}$ , nearby solutions rotate around it in a clockwise manner. This “rotation” in the  $\theta-v$  coordinates corresponds to a flight path which wobbles: the nose angle alternately increases and decreases, as does the velocity. For  $v$  large, it seems that the solutions move from left to right— $\theta$  constantly increases, which means that the glider is always increasing its nose angle. This means it is doing loops.

We can also use the `DEplot` command from Maple’s `DEtools` package to produce a direction field. The direction field is obtained from the vector field by rescaling all nonzero vectors to be of unit length. Sometimes discarding the extra information makes it easier to see what is happening. In addition, `DEplot` has other very useful options, and will be one of our primary tools throughout this chapter.

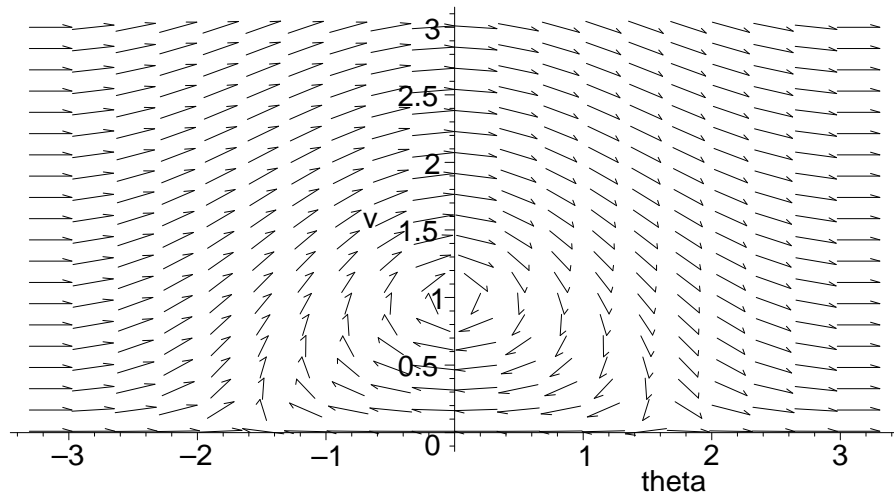
For convenience, we will name our system of differential equations as `phug`, so that we needn’t keep retyping them.

```
> phug := [ diff(theta(t),t) = (v(t)^2 - cos(theta(t)))/v(t),
            diff(v(t),t)      = -sin(theta(t)) - R*v(t)^2 ];
R:=0;
```

$$phug := \left[ \frac{\partial}{\partial t} \theta(t) = \frac{v(t)^2 - \cos(\theta(t))}{v(t)}, \frac{\partial}{\partial t} v(t) = -\sin(\theta(t)) \right]$$

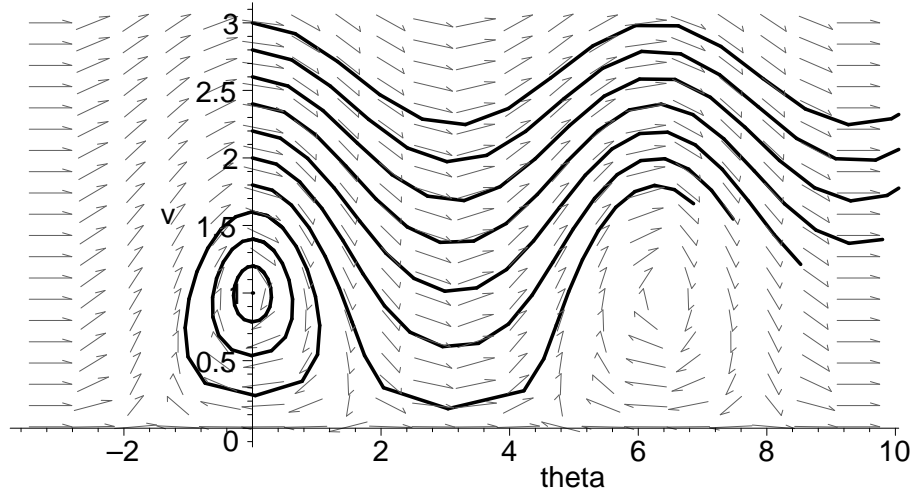
$$R := 0$$

```
> with(DEtools):
DEplot( phug, [theta(t), v(t)], t=0..1, theta=-Pi..Pi, v=0.01..3);
```



We can also ask maple to plot a number of solutions corresponding to several initial conditions. Rather than typing in 11 initial conditions, we can use `seq` to generate them. We tell maple that the solution curves should be black, because otherwise it uses a sickly yellow. If we had wanted each solution curve to be a different color, we could use something like `linecolor=[seq(COLOR(HUE,i/10),i=0..10)]`.

```
> DEplot( phug, [theta(t), v(t)], t=0..5,
  [seq([theta(0)=0, v(0)=1+0.2*i], i=0..10)],
  theta=-Pi..3*Pi, v=0.01..3, linecolor=black );
```



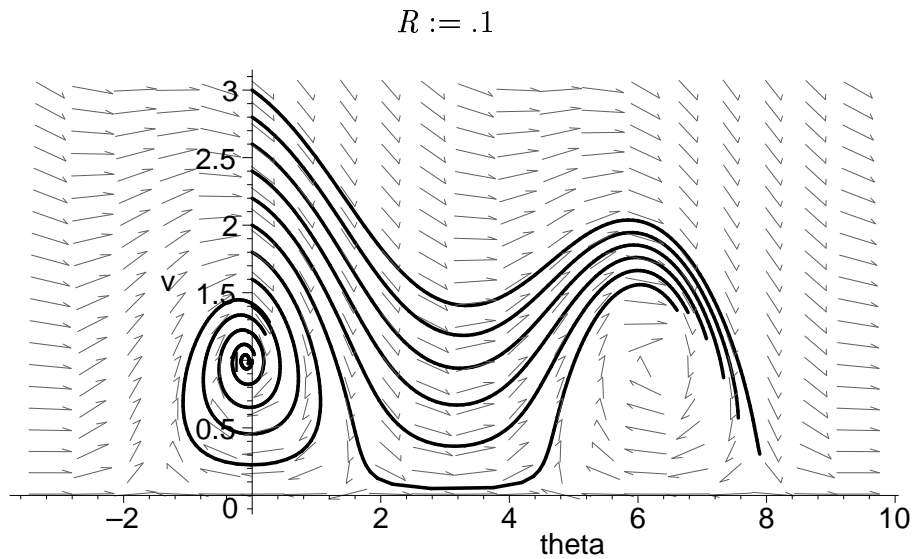
The above figure confirms our earlier impressions about the “wobbling” and “looping” solutions. We’ll see how to get Maple to demonstrate this more explicitly in section 5, where we will plot the path of the glider through space (i.e., in  $x$ - $y$  coordinates).

Just to see how things differ if we make  $R$  nonzero, let’s change  $R$  and get Maple to show us what happens. We need to adjust the stepsize to get a reasonable accuracy; what precisely this does is described in the next section.

```

> R:=0.1;
DEplot( phug, [theta(t), v(t)], t=0..5,
        [seq([theta(0)=0, v(0)=1+0.2*i], i=0..10)],
        theta=-Pi..3*Pi, v=0.01..3, linecolor=black, stepsize=0.1 );

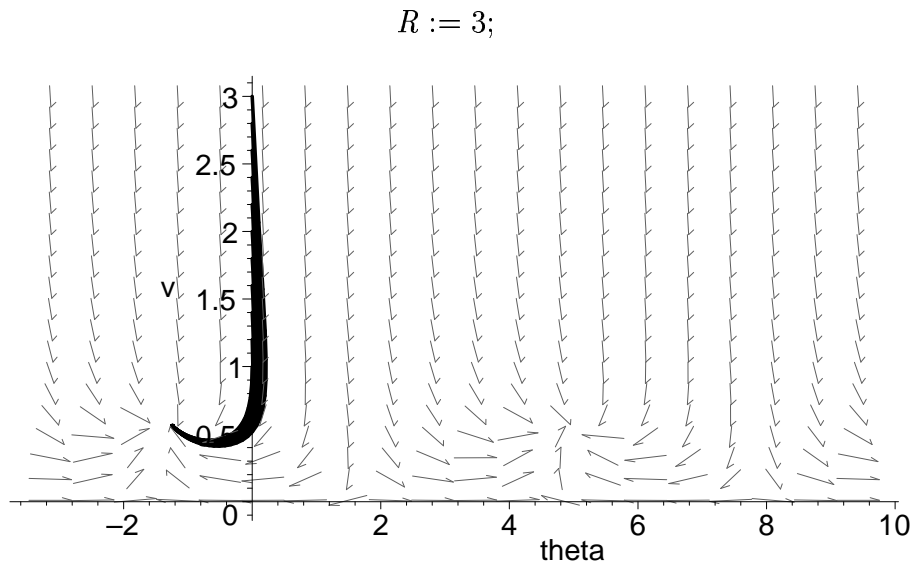
```



```

> R:=3;
DEplot( phug, [theta(t), v(t)], t=0..5,
        [seq([theta(0)=0, v(0)=1+0.2*i], i=0..10)],
        theta=-Pi..3*Pi, v=0.01..3, linecolor=black, stepsize=0.1 );

```



The behavior of the glider appears to be dramatically affected by the drag parameter  $R$ . One of our main goals in this chapter will be to classify the types of solutions that are possible as a function of the drag parameter  $R$ .

### 3 Existence of Solutions

We've been acting as though just by specifying an initial condition, there must be a solution, and it must be unique (that is, the only one corresponding to that initial condition). And, in fact, this is typically true for any "nice" differential equation. But which differential equations are "nice" enough?

We won't prove the theorem here, but we will state it. The reader is encouraged to look up the proof in any differential equations text, such as [BDH] or [HW].

**Theorem 3.1** *Consider the differential equation*

$$\frac{d\vec{x}}{dt} = \vec{F}(\vec{x}, t).$$

*If  $\vec{F}$  is continuous on an open set containing  $(\vec{x}_0, t_0)$ , then there is a solution  $\vec{x}(t)$  which is defined on some interval  $(t_0 - \epsilon, t_0 + \epsilon)$  which satisfies the initial value problem*

$$\frac{d\vec{x}}{dt} = \vec{F}(\vec{x}, t) \quad \vec{x}(t_0) = \vec{x}_0.$$

*Furthermore, if  $\vec{F}$  has continuous partial derivatives, this solution is unique.*

In fact, the hypothesis can be weakened a little bit and still preserve the uniqueness. As long as  $\vec{F}$  is at least Lipschitz in  $\vec{x}$ , the solution will be unique. A function  $\vec{F}$  is called Lipschitz if there is some  $K$  so that

$$|\vec{F}(\vec{x}_1, t) - \vec{F}(\vec{x}_2, t)| < K|\vec{x}_1 - \vec{x}_2|$$

for all  $t$ ,  $\vec{x}_1$ , and  $\vec{x}_2$ . This condition says, roughly, that  $\vec{F}$  doesn't spread out too quickly. All functions with continuous partials are automatically Lipschitz.

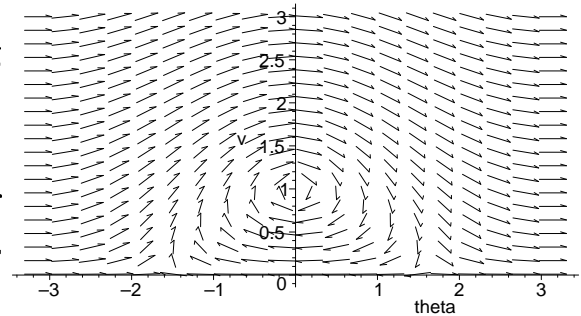
This means that for us, as long as we stay away from the place where our differential equation isn't defined (i.e. ensure that  $v > 0$ ), we can be assured that there is a solution through every point, and that solution is unique. Notice that since our equations are autonomous (the right-hand side has no explicit dependence on  $t$ ), uniqueness of solutions means that solutions cannot cross in the  $(\theta, v)$ -plane.

### 4 Numerical Methods

What does Maple do when you ask it to display a solution of a differential equation with DEplot? Obviously, it cannot solve the equations analytically, since solutions don't always exist in closed form. Instead, it approximates them numerically.

## 4.1 Euler's method

To get an idea of how this can be done, take a look again at the direction field for the glider. By looking at how the arrows point, you can almost guess how the solution must turn. From a given initial condition, move in the direction of the arrow to the next position. At that new point, step in the direction of the vector at that point, and so on. This is the idea behind the simplest numerical integration scheme, called Euler's method.



More precisely, given a differential equation and an initial condition

$$\frac{d\vec{x}}{dt} = \vec{F}(\vec{x}, t) \quad \varphi(t_0) = \vec{x}_0,$$

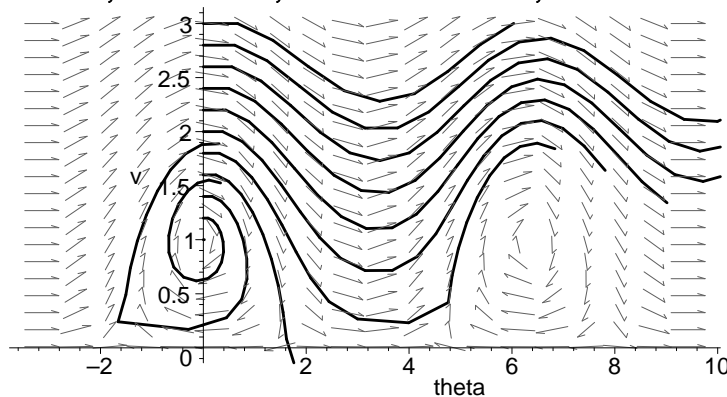
the goal is to construct a sequence of points  $\vec{x}_0, \vec{x}_1, \vec{x}_2, \dots, \vec{x}_n$  which approximate the solution  $\varphi$  at times  $t_0, t_1, t_2, \dots, t_n$ . We take our initial point to be the initial condition, and choose some small number  $h > 0$  called the *stepsize*. Then we set

$$\vec{x}_{i+1} = \vec{x}_i + h\vec{F}(\vec{x}_i, t_i) \quad t_{i+1} = t_i + h$$

for  $0 \leq i \leq n - 1$ . This very simple scheme works best for small values of  $h$ ; the error in the approximation for a fixed time interval is proportional to  $h$ . This means that to improve the accuracy of a numerical solution by a factor of 10, we need to do about 10 times more work.

Maple doesn't typically use Euler's method because the errors accumulate too quickly. To see this, we reproduce one of the earlier figures from §2 using Euler's method.

```
> DEplot( phug, [theta(t), v(t)], t=0..5,
  [seq([theta(0)=0, v(0)=1+0.2*i], i=0..10)],
  theta=-Pi..3*Pi, v=0.01..3, linecolor=black, method=classical[foreuler]);
```



While the errors in the above figures are considerable, using a much smaller stepsize (say, `stepsize=0.01`) will essentially reproduce the earlier figure, with a lot more computation.



## 4.2 Numerical Solutions, Numerical Integration, and Runge-Kutta

To get some ideas about improving on Euler's method, let's first notice that Euler's method corresponds exactly to the left-hand Riemann sum when computing a numerical approximation to an integral. You may recall from calculus that the left-hand sum converges to the value of the integral, but you must take a large number of subdivisions to get any accuracy.

A more efficient method is the trapezoid rule, which is the average of the left-hand and right-hand sum. There is a numerical scheme for integrating ODEs which corresponds to this, known variously as the *improved Euler method* or the *Heun formula*. In the trapezoid rule, we compute the function at the left and right endpoints of the function and average the result; the analog of this would require us to average the value of the vector field at our current point and at the next point. Perhaps you see the circularity implicit in that reasoning: to know the next point, we would need to solve the differential equation. So, instead we fudge a bit: we get our approximation of the next point by using Euler's method. More specifically, using the notation of §4.1, given an approximation  $\vec{x}_i$  at time  $t_i$ , we find the next one using the scheme below.

$$\begin{array}{ll}
 t_{i+1} &= t_i + h \\
 \vec{k}_\ell &= \vec{F}(\vec{x}_i, t_i) && \text{slope at the left side of the interval} \\
 \vec{k}_r &= \vec{F}(\vec{x}_i + h\vec{k}_\ell, t_{i+1}) && \text{slope at the right side of the interval} \\
 \vec{x}_{i+1} &= \vec{x}_i + h(\vec{k}_\ell + \vec{k}_r)/2
 \end{array}$$

To do the improved Euler method, we need to evaluate the vector field twice as often as the regular Euler method, but there is a big gain in accuracy. The error of the approximation for a fixed time interval is proportional to  $h^2$ . Thus, if we decrease  $h$  by a factor of 10, we should expect to reduce the error 100-fold. You can get Maple to use the improved Euler method by specifying `method=classical[heunform]` as an argument to `DEplot` and related commands.

You might remember from your calculus course that Simpson's rule can get very accurate numerical approximations for integrals with a very few number of points. This is because Simpson's rule is a fourth-order method: the error in the approximation is proportional to  $h^4$ . For Simpson's rule, we evaluate the function at the right endpoint, the left endpoint, and the middle, and then take a weighted average of those values.

The analog of Simpson's rule for differential equations is called the Runge-Kutta method<sup>2</sup>, developed by the German mathematicians C. Runge and W. Kutta at the end of the nineteenth

---

<sup>2</sup>In fact, there is a whole family of such methods, of various orders. But generally, if someone says "Runge-Kutta" without specifying the order, they mean this fourth-order Runge-Kutta.

century. It takes a linear combination of four slopes, as below:

$$\begin{array}{ll}
 t_{i+1} &= t_i + h \\
 \vec{k}_\ell &= \vec{F}(\vec{x}_i, t_i) && \text{slope at the left side of the interval} \\
 \vec{k}_m &= \vec{F}(\vec{x}_i + \frac{h}{2}\vec{k}_\ell, t_i + \frac{h}{2}) && \text{slope at the middle of the interval} \\
 \vec{k}_c &= \vec{F}(\vec{x}_i + \frac{h}{2}\vec{k}_m, t_i + \frac{h}{2}) && \text{corrected slope at the middle of the interval} \\
 \vec{k}_r &= \vec{F}(\vec{x}_i + h\vec{k}_c, t_{i+1}) && \text{slope at the right side of the interval} \\
 \vec{x}_{i+1} &= x_i + h(\vec{k}_\ell + 2\vec{k}_m + 2\vec{k}_c + \vec{k}_r)/6
 \end{array}$$

When the slope depends only on  $t$  and not on  $\vec{x}$ , then  $\vec{k}_m$  and  $\vec{k}_c$  are the same, and the formula reduces exactly to that of Simpson's rule.

Runge-Kutta is probably the most commonly used method of numerical integration, because of its generally high accuracy. Like Simpson's rule, it is a fourth-order method, so the error is proportional to  $h^4$ . Maple has several numerical methods for ODEs built in to it; see the help page on `dsolve[numeric]` for more information about them; the ones we have described are "classical" methods, and are described (along with others) on Maple's help page for `dsolve[classical]`. Unless asked to do otherwise, Maple's `DEplot` command uses the fourth-order Runge-Kutta method described above (Maple's name for this method is `classical[rk4]`), although the `numeric` option of `dsolve` defaults to a variation called *fourth-fifth order Runge-Kutta-Fehlberg*, which Maple refers to as `rkf45`.

## 5 Seeing the flight path

In §2, we exhibited a number of solutions to the phugoid model with various initial conditions and values for the drag, as plots of the parametric curves  $[\theta(t), v(t)]$ . From these, we could deduce an approximate path of the glider through space via reasoning like " $\theta(t)$  is monotonically increasing, so the glider must be doing loop-the-loops". We will now discuss how to solve for the position  $[x(t), y(t)]$  directly.

First, notice that if we know the value of  $v$  and  $\theta$  at some time  $t$ , we also have  $\frac{dx}{dt}$  and  $\frac{dy}{dt}$  at that  $t$ , since

$$\frac{dx}{dt} = v \cos \theta \quad \text{and} \quad \frac{dy}{dt} = v \sin \theta.$$

So, one way to find  $x(t)$  and  $y(t)$  would be to solve the original ODE (numerically) for  $\theta(t)$  and  $v(t)$ , obtain a new differential equation, and then solve that numerically for  $x(t)$  and  $y(t)$ .

Aside from the fact that there is an easier way, there is a serious practical problem with this approach. When we get a numerical solution  $[\theta(t), v(t)]$ , this is a collection of points  $(\theta_0, v_0, t_0), (\theta_1, v_1, t_1), \dots, (\theta_n, v_n, t_n)$  along a curve. This means we only know the vector field for  $(x, y)$  at those points. It is highly unlikely that these will be the points we need to use Runge-Kutta. We could use these to approximate our  $(x, y)$ -solution using Euler's method,

although even then we would have to program it ourselves rather than use Maple's built-in method.

If you think for just a moment, you should realize that there is no reason we must find  $\theta(t)$  and  $v(t)$  before attempting to find  $x(t)$  and  $y(t)$ . Instead, we can augment the original system, adding in the new equations for  $\dot{x}$  and  $\dot{y}$ . Then, we can use numerical methods to solve for  $[\theta(t), v(t), x(t), y(t)]$  all in one step. The methods we discussed in §4.2 apply equally well to vector fields in any number of variables.

Consequently, we rewrite our phugoid system as

$$\frac{d\theta}{dt} = \frac{v^2 - \cos\theta}{v} \quad \frac{dv}{dt} = -\sin\theta - Rv^2 \quad \frac{dx}{dt} = v \cos\theta \quad \frac{dy}{dt} = v \sin\theta.$$

Using this system, we need to specify our initial conditions as  $(\theta_0, v_0, x_0, y_0)$ , and each solution is a curve in  $\mathbb{R}^4$  parameterized by time. We can take a projection of this curve to get our solutions in  $\theta$ - $v$  or  $x$ - $y$  coordinates.

Asking Maple to solve this new system is almost the same as before, but we must specify the additional variables and use the `scene` option to `DEplot` to tell Maple which projection we want to see. Note that we reset  $R$  to ensure that it appears in the equations as an arbitrary constant, rather than with any previously assigned value.

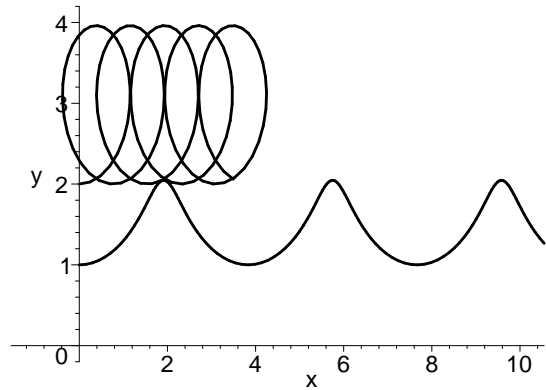
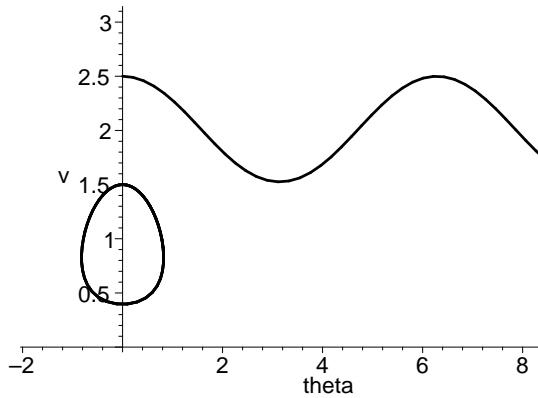
```
> R:=‘R‘;
xphug :=[ diff(theta(t),t) = (v(t)^2 - cos(theta(t)))/v(t),
          diff(v(t),t) = -sin(theta(t)) - R*v(t)^2 ,
          diff(x(t),t) = v(t)*cos(theta(t)),
          diff(y(t),t) = v(t)*sin(theta(t))];

xphug := [  $\frac{\partial}{\partial t} \theta(t) = \frac{v(t)^2 - \cos(\theta(t))}{v(t)}$ ,  $\frac{\partial}{\partial t} v(t) = -\sin(\theta(t)) - Rv(t)^2$ ,
           $\frac{\partial}{\partial t} x(t) = v(t) \cos(\theta(t))$ ,  $\frac{\partial}{\partial t} y(t) = v(t) \sin(\theta(t))$ ]
```

Now we ask Maple to plot two solution curves in both  $\theta, v$  coordinates and  $x, y$ . One is started with an initial angle of  $\theta = 0$  and velocity  $v = 1.5$  at a height  $y(0) = 1$ ; the other starts with  $\theta = 0$ ,  $v = 2.5$ , and initial height  $y = 2$ . We also use `display` with an `array` to show the two scenes side by side.

```
> R:=0;
plots[display](array([
  DEplot(xphug, [theta,v,x,y], t=0..15,
    [[theta(0)=0,v(0)=1.5,x(0)=0,y(0)=1], [theta(0)=0,v(0)=2.5,x(0)=0,y(0)=2]],
    theta=-Pi..5*Pi/2, v=0.1..3, x=-1..10, y=0..4,
    linecolor=black, stepsize=0.1, scene=[theta,v]),

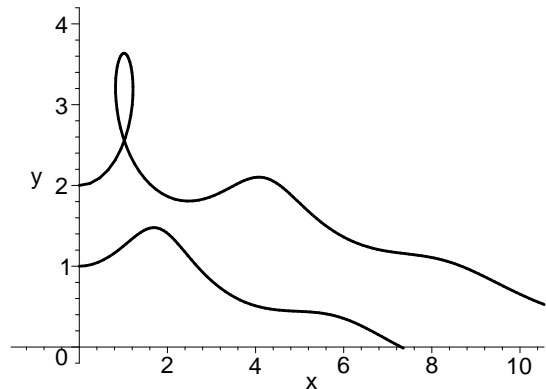
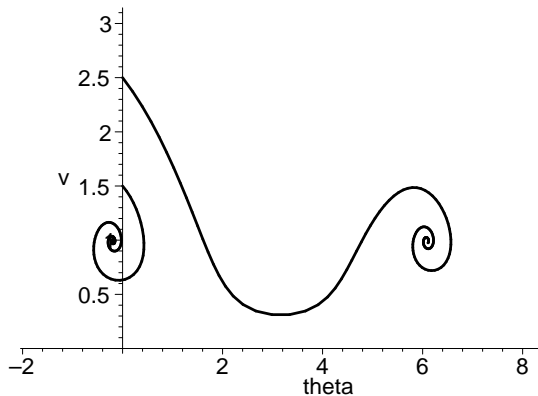
  DEplot(xphug, [theta,v,x,y], t=0..15,
    [[theta(0)=0,v(0)=1.5,x(0)=0,y(0)=1], [theta(0)=0,v(0)=2.5,x(0)=0,y(0)=2]],
    theta=-Pi..5*Pi/2, v=0.1..3, x=-1..10, y=0..4,
    linecolor=black, stepsize=0.1, scene=[x,y]),
  ]));
```

$R := 0;$ 


If we change the value of  $R$  to 0.2 and repeat the command, we can see how things change when there is friction. Notice how spiraling towards the fixed point in the  $\theta$ - $v$  coordinates corresponds to a plane which is diving with an oscillatory motion.

```
> R:=0.2;
plots[display](array([
  DEplot(xphug, [theta,v,x,y], t=0..15,
    [[theta(0)=0,v(0)=1.5,x(0)=0,y(0)=1],[theta(0)=0,v(0)=2.5,x(0)=0,y(0)=2]],
    theta=-Pi..5*Pi/2, v=0.1..3, x=-1..10, y=0..4,
    linecolor=black, stepsize=0.1, scene=[theta,v]),

  DEplot(xphug, [theta,v,x,y], t=0..15,
    [[theta(0)=0,v(0)=1.5,x(0)=0,y(0)=1],[theta(0)=0,v(0)=2.5,x(0)=0,y(0)=2]],
    theta=-Pi..5*Pi/2, v=0.1..3, x=-1..10, y=0..4,
    linecolor=black, stepsize=0.1, scene=[x,y]),
)]));
```

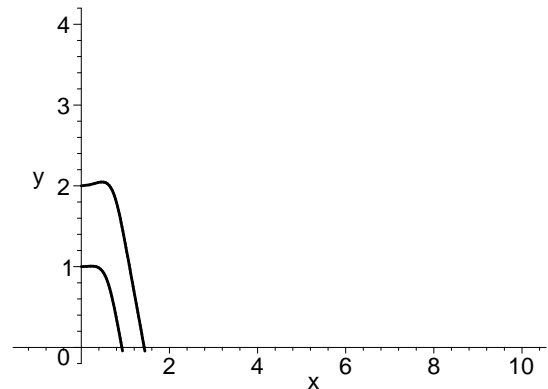
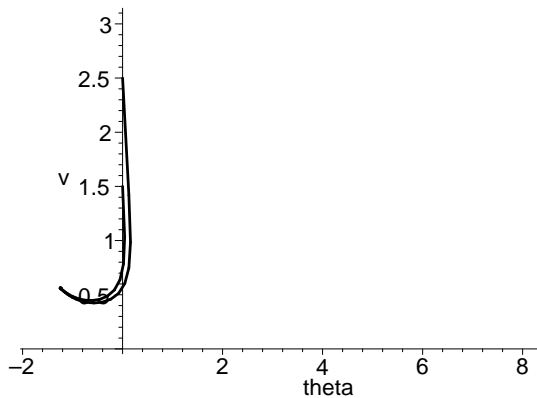
 $R := 0.2;$ 


Finally, we redo the pictures with a large amount of friction ( $R = 3$ ). Note how the glider acts more like a rock than a glider, and there is no oscillation at all.

```
> R:=3;
plots[display](array([
  DEplot(xphug, [theta,v,x,y], t=0..15,
    [[theta(0)=0,v(0)=1.5,x(0)=0,y(0)=1],[theta(0)=0,v(0)=2.5,x(0)=0,y(0)=2]],
    theta=-Pi..5*Pi/2, v=0.1..3, x=-1..10, y=0..4,
    linecolor=black, stepsize=0.1, scene=[theta,v]),

  DEplot(xphug, [theta,v,x,y], t=0..15,
    [[theta(0)=0,v(0)=1.5,x(0)=0,y(0)=1],[theta(0)=0,v(0)=2.5,x(0)=0,y(0)=2]],
    theta=-Pi..5*Pi/2, v=0.1..3, x=-1..10, y=0..4,
    linecolor=black, stepsize=0.1, scene=[x,y]),
  ]));
```

$R := 3;$



## 6 Fixed Point Analysis

So far, the cases we've looked at seem to have a number of common features. Each has a solution with  $\theta(t)$  and  $v(t)$  constant; there are solutions which immediately tend toward the constant solution (or oscillate around it), and solutions which do some number of loops first (for  $R = 3$ , we didn't demonstrate this— if the initial angle is 0, a solution needs an initial velocity larger than 86.3 before it will do a loop). How do the possible solutions depend on the value of  $R$ ? For example, the behavior of solutions for  $R = 0.1$  and  $R = 0.2$  are quite similar, but are dramatically different from  $R = 0$  and  $R = 3$ .

One way we can get a better handle on exactly how the behavior of the solutions depends on  $R$  is to examine what happens to solutions near the constant solution. If we look in the  $\theta$ - $v$  plane, this solution corresponds to a single point, which is often referred to as a “fixed point”. We explicitly found the fixed point  $\{\theta(t) = 0, v(t) = 1\}$  solution for  $R = 0$  in §2 by noticing

that whenever  $\dot{\theta} = 0$  and  $\dot{v} = 0$ , we must have a constant solution (and vice-versa). Using `solve`, we can get maple to tell us the where the fixed point is for arbitrary values of  $R$ . We use `convert` to insist that maple represent the result as a radical, instead of using the `RootOf` notation.

```
> R:='R':
  FixPoint:=convert(
    solve({(v^2-cos(theta))/v=0, -sin(theta)-R*v^2=0}, {theta,v}),
    radical);
```

$$FixPoint := \{v = \left(\frac{1}{R^2 + 1}\right)^{(1/4)}, \theta = \arctan\left(-R \sqrt{\frac{1}{R^2 + 1}}, \sqrt{\frac{1}{R^2 + 1}}\right)\}$$

So, we see that there is a fixed point for all values of  $R$  (since  $1 + R^2$  is always positive, the radical always takes real values). For  $R > 0$ ,  $\theta(t)$  at this fixed solution is negative, so this corresponds to a diving solution. As  $R$  increases, the angle of the dive becomes steeper and steeper.

What can we say about behaviour of solutions *near* this fixed solution?

In order to answer this, notice that we can apply Taylor's theorem to the right-hand side of our system of differential equations. That is, think of our differential equation as being in the vector form

$$\frac{d\vec{X}}{dt} = \vec{F}(\vec{X}),$$

where (in our case)  $\vec{X}$  is  $[\theta(t), v(t)]$ , and  $\vec{F}(X) = [(v^2 - \cos(\theta))/v, -\sin(\theta) - Rv^2]$ . Taylor's theorem says that for  $\vec{X}$  near  $\vec{X}_0$ ,

$$\vec{F}(\vec{X}) = \vec{F}(\vec{X}_0) + D\vec{F}(\vec{X}_0)(\vec{X} - \vec{X}_0) + \text{higher order terms},$$

where  $D\vec{F}(\vec{X}_0)$  is the Jacobian of  $\vec{F}$  evaluated at  $\vec{X}_0$ . Since  $\vec{X}_0$  is a fixed point,  $\vec{F}(\vec{X}_0)$  is the zero vector. Hence, we can get a good idea of what is happening close to the fixed point by studying the *linear* system of differential equations given by the derivative at the fixed point.

First, we give a brief refresher about linear  $2 \times 2$  systems of equations. The reader is referred to the texts in the references for more comprehensive treatment of this topic.

## 6.1 Linear Systems of ODEs

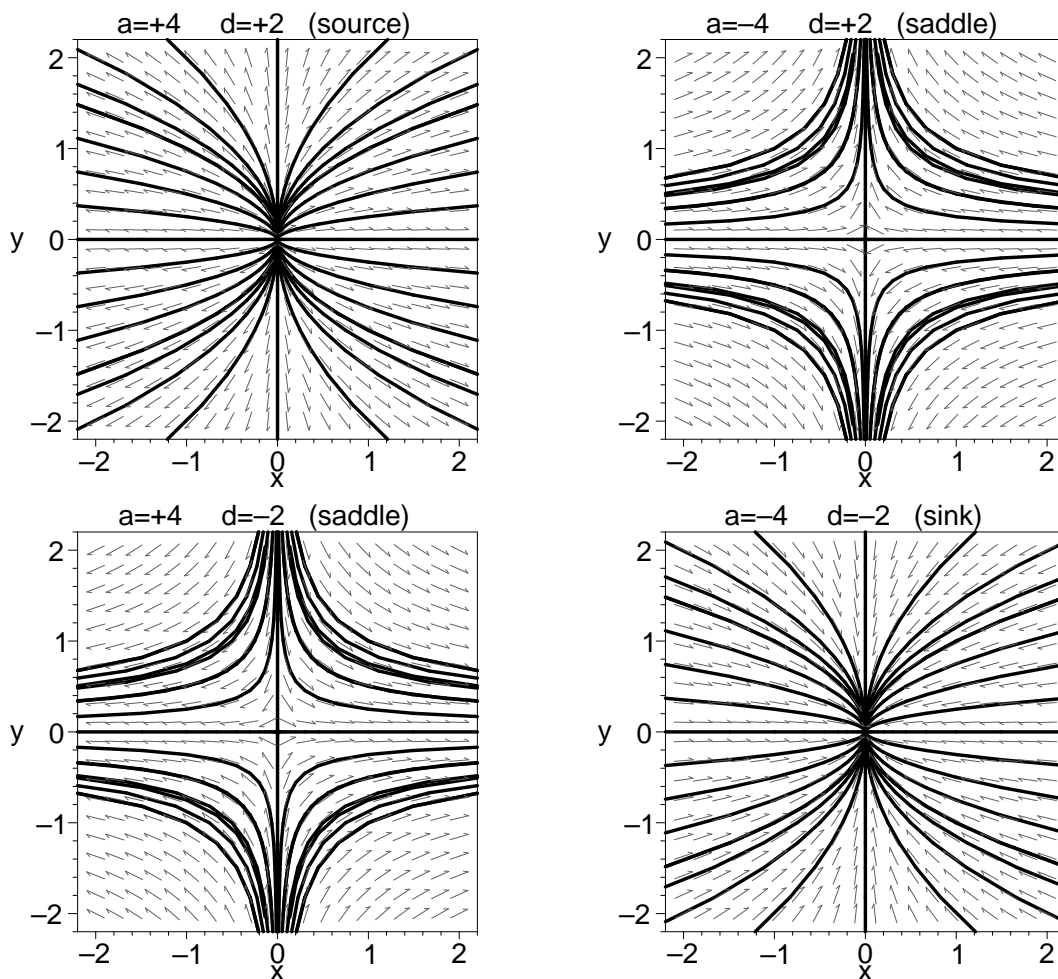
Linear differential equations (with constant coefficients) are the simplest ones to solve and understand. Such equations have the form  $\frac{dX}{dt} = AX$ , where  $A$  is an  $n \times n$  matrix and  $X$  is an  $n$ -vector. We'll be content with  $n = 2$ .

First, notice that  $X = 0$  is always a fixed point for this system, and it is the only one. Secondly, if  $X_1(t)$  and  $X_2(t)$  are solutions to the system, then so is  $X_1(t) + X_2(t)$ . This property is called *superposition* and will be very useful.

Let's look at the simplest case for the  $2 \times 2$  system:  $A = \begin{pmatrix} a & 0 \\ 0 & d \end{pmatrix}$ , that is,

$$\begin{aligned}\dot{x} &= ax \\ \dot{y} &= dy\end{aligned}$$

If we have an initial condition with  $y(0) = 0$ , then  $y(t) = 0$ . This means that the problem reduces to the one-dimensional case, and  $x(t) = x(0)e^{at}$ . If  $a > 0$ , the solution moves away from the origin along the  $x$ -axis as  $t$  increases; if  $a < 0$ , the solution moves toward the origin. If the initial condition is on the  $y$ -axis, the sign of  $d$  controls whether the solution moves away ( $d > 0$ ) or towards ( $d < 0$ ) the origin. Because of the superposition property mentioned before, an initial condition with  $x(0) \neq 0$  and  $y(0) \neq 0$  gives rise to a solution which is a combination of these two behaviours.



A large number of linear systems behave very similarly to the cases above, except the straight-line solutions may not lie exactly along the coordinate axes.

Suppose there is a solution of  $\dot{X} = AX$  which lies along a straight line, that is, along some vector  $\vec{v}$ . Because the system is linear, this means the tangent vectors to this solution must be of the form  $\lambda X$  for some number  $\lambda$ . Such a number  $\lambda$  is called an *eigenvalue*, and the corresponding vector  $\vec{v}$  is the corresponding *eigenvector*. Note that in the case above, we had eigenvalues  $a$  and  $d$  with eigenvectors  $[1, 0]$  and  $[0, 1]$ .

To find the eigenvalues, we need to solve

$$AX = \lambda X, \quad \text{or equivalently,} \quad (A - \lambda I)X = 0.$$

This can only happen if  $X$  is the zero vector, or the determinant of  $A - \lambda I$  is zero. In the latter case, we must have

$$(a - \lambda)(d - \lambda) - bc = 0,$$

That is,

$$\lambda^2 - (a + d)\lambda + (ad - bc) = 0.$$

Using the quadratic formula, we see that the eigenvalues must be

$$\lambda = \frac{(a + d) \pm \sqrt{(a + d)^2 - 4(ad - bc)}}{2}.$$

The quantity  $a + d$  is called the *trace* of  $A$  (more generally, the trace of a matrix is the sum of diagonal entries), and  $ad - bc$  is the determinant of  $A$ . The eigenvalues of a  $2 \times 2$  matrix can be expressed in terms of the trace and the determinant<sup>3</sup> as

$$\frac{\text{Tr}A \pm \sqrt{(\text{Tr}A)^2 - 4 \det A}}{2}.$$

We'll use this form again in a little while.

For linear systems, the eigenvalues determine the ultimate fate of solutions to the ODE. From the above, it should be clear that there are always two eigenvalues for a  $2 \times 2$  matrix (although sometimes there might be a double eigenvalue, when  $(\text{Tr}A)^2 - 4 \det A = 0$ ). Let's call them  $\lambda_1$  and  $\lambda_2$ .

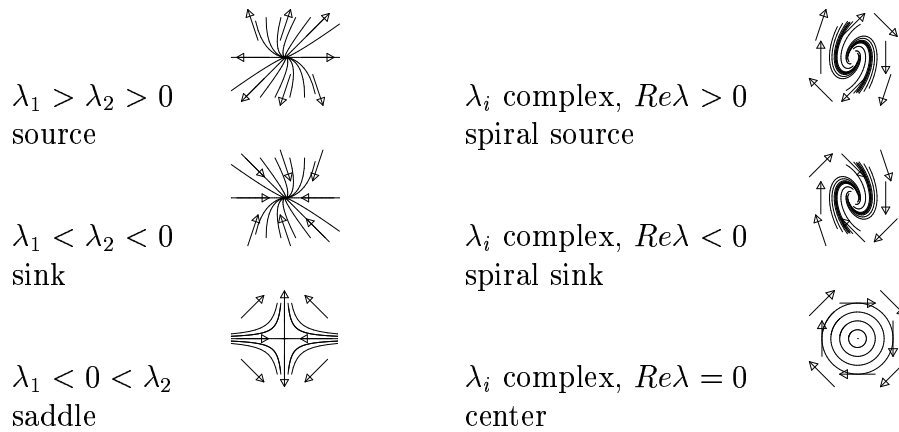
We've already seen prototypical examples where the eigenvalues are real, nonzero, and distinct. But what if the eigenvalues are complex conjugate (which happens when  $(\text{Tr}A)^2 < 4 \det A$ )? In this case, there is no straight line solution in the reals. Instead, solutions turn around the origin. If the real part of the eigenvalues is positive, solutions spiral away from the origin, and if it is negative, they spiral towards it. If the eigenvalues are purely imaginary, the solutions neither move in nor out; rather, they circle around the origin.

We can summarize this in the table below. We are skipping over dealing with the degenerate cases, when  $\lambda_1 = \lambda_2$  and when one of the eigenvalues is zero.

---

<sup>3</sup>It is true for all matrices that the trace is the sum of the eigenvalues, and the determinant is their product.

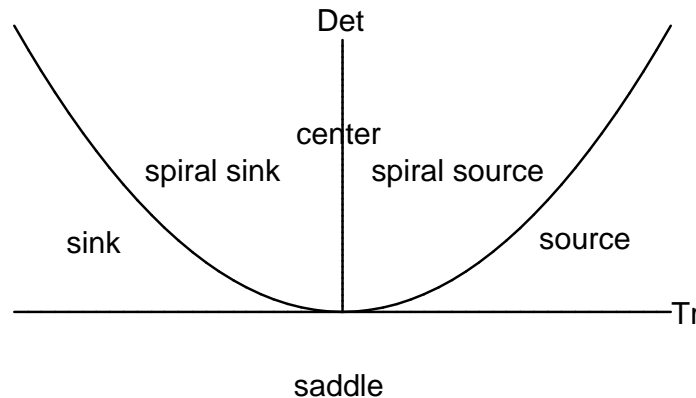




Finally, we remark that there is a convenient way to organize this information into a single diagram. Since the eigenvalues of  $A$  can be written as

$$\frac{TrA \pm \sqrt{(TrA)^2 - 4 detA}}{2},$$

we can consider the matrix  $A$  as being a point in the Trace-Determinant plane. Then the curve where  $(TrA)^2 = 4 detA$  is a parabola. If  $A$  is above this parabola, it has complex eigenvalues. Furthermore, if the determinant is negative, we must have a saddle because<sup>4</sup>  $\sqrt{(TrA)^2 - 4 detA} > TrA$ , which means there is one positive and one negative eigenvalue. Finally, if the determinant is positive, the eigenvalues (or their real part, if they are complex) is the same sign as the trace of  $A$ . We summarize this in the diagram below.



<sup>4</sup>Since the determinant is the product of the eigenvalues, the only way it can be negative is if they are of opposite signs.

## 6.2 Fixed Points for the Glider

We now return to our nonlinear system, and look at the linearization near the fixed point. As we saw earlier in this section, for every value of  $R$  the Phugoid model

$$\frac{d\theta}{dt} = \frac{v^2 - \cos \theta}{v} \quad \frac{dv}{dt} = -\sin \theta - Rv^2$$

has a fixed point at

$$\theta = -\arctan\left(R\sqrt{\frac{1}{1+R^2}}\right) \quad v = \sqrt[4]{\frac{1}{1+R^2}}.$$

We now calculate the Jacobian matrix, using Maple.

```
> with(linalg):
R:='R':
J:=jacobian( [(v^2-cos(theta))/v, -sin(theta)-R*v^2], [theta, v]);
```

$$J := \begin{bmatrix} \frac{\sin(\theta)}{v} & 2 - \frac{v^2 - \cos(\theta)}{v^2} \\ -\cos(\theta) & -2Rv \end{bmatrix}$$

Once we know the trace and determinant of the Jacobian matrix at the fixed point, we can use that information to determine how its type (sink, source, saddle, etc.) depends on  $R$ . Rather than solve for the eigenvalues directly, we will compute the trace and determinant, which have a simpler form in this case.

First, let's find the fixed point again, issuing the same command we did earlier. Then, we'll substitute that into the result of `trace` and `det` to calculate the trace and determinant at the fixed point.

```
> FixPoint:=convert(
    solve({(v^2-cos(theta))/v=0, -sin(theta)-R*v^2=0}, {theta,v}),
    radical);
```

$$\text{FixPoint} := \left\{ \theta = \arctan\left(-R\sqrt{\frac{1}{R^2+1}}, \sqrt{\frac{1}{R^2+1}}\right), v = \left(\frac{1}{R^2+1}\right)^{(1/4)} \right\}$$

```
> trfix := simplify(subs(FixPoint, trace(J)));
detfix := simplify(subs(FixPoint, det(J)));
```

$$\text{trfix} := \frac{-3R}{(R^2+1)^{1/4}}$$

$$\text{detfix} := 2\sqrt{R^2+1}$$

So, we see that for all real values of  $R$ , the determinant is positive (meaning saddles are impossible). Furthermore, for  $R > 0$ , the trace is negative, so the fixed point is always a sink. (When  $R = 0$ , the trace is zero, which means it is possible for this fixed point to be a center.<sup>5</sup>) From the plots we made earlier, we expect that for small values of  $R$ , we will have a spiral sink, but it will probably not spiral for sufficiently large values of  $R$ . We can check this by solving for the point where  $Tr(J)^2 = 4\det(J)$ , which happens when  $R = 2\sqrt{2}$ .

```
> solve( (trfix)^2 = 4*detfix, R);
```

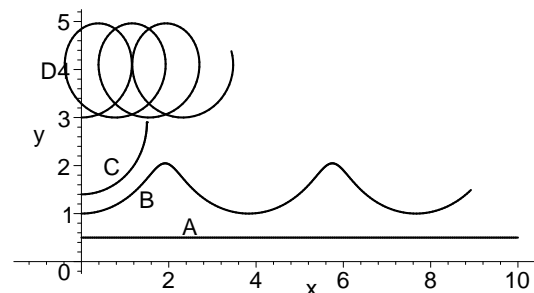
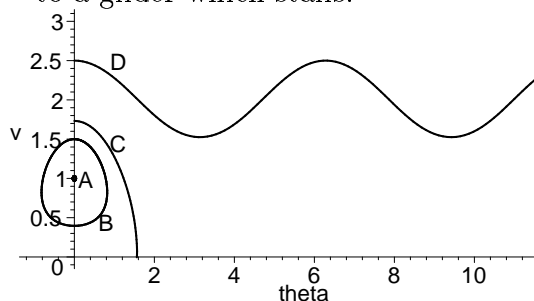
$$2\sqrt{2}, -2\sqrt{2}$$

## 7 Qualitative Classification of Solutions

This analysis of the previous section allows us to completely classify what happens in the Phugoid model for all  $R \in [0, \infty)$ . In all cases, the fixed point in  $(\theta, v)$  coordinates corresponds to the solution

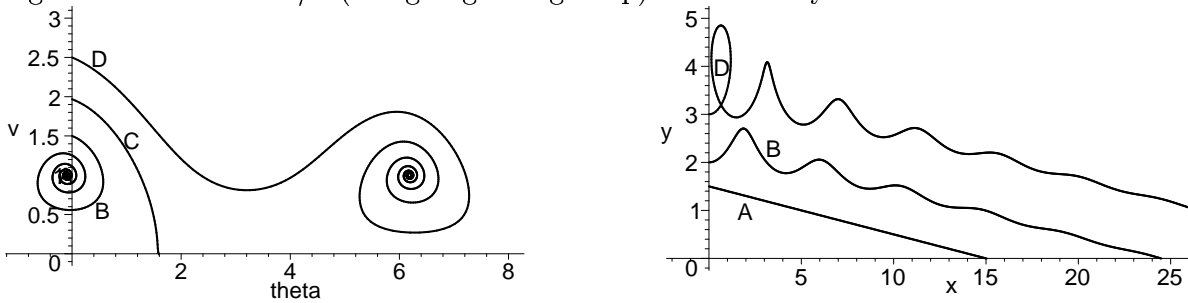
$$\theta(t) = -\arctan\left(R\sqrt{\frac{1}{1+R^2}}\right) \quad v(t) = \sqrt[4]{\frac{1}{1+R^2}} \quad x(t) = \frac{t}{(1+R^2)^{\frac{3}{4}}} \quad y(t) = \frac{-Rt}{(1+R^2)^{\frac{3}{4}}}$$

$R = 0$ : In this case, there is no drag. The fixed point at  $\theta = 0, v = 1$  corresponds to a glider flying level with a constant speed (labeled A in the figures below). The fixed point is a center in  $\theta, v$  coordinates: nearby solutions are closed curves and correspond to a glider with an oscillatory path, alternately diving and climbing (labeled B). For initial conditions further away from the fixed point,  $v(t)$  oscillates, but  $\theta(t)$  constantly increases (see D). Such solutions correspond to a glider endlessly looping, and a pilot with a severe case of nausea. Between these two is a solution which cannot be continued beyond a certain time (C), because  $v(t)$  becomes zero and our equations are no longer defined. This corresponds to a glider which stalls.

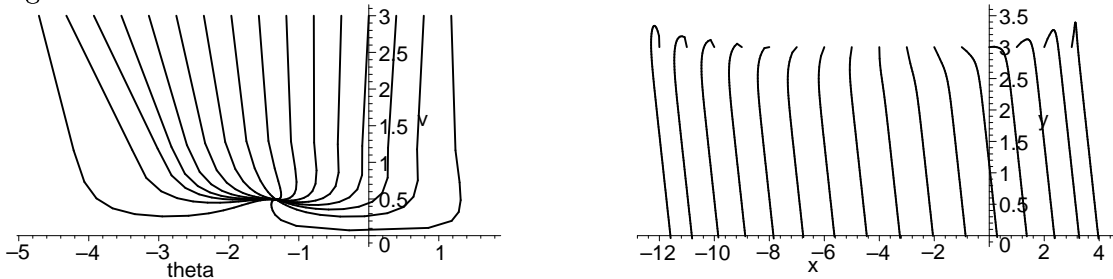


<sup>5</sup>We didn't address this point. If the trace is nonzero, it is impossible to have a center. But if the trace is zero, then the higher-order terms in the Taylor series play a role, and anything is possible. In fact, it can be shown that the fixed point for  $R = 0$  is indeed a center. One way to do this is to show that the quantity  $E(\theta, v) = v^3 - 3v \cos \theta$  must be constant for any solution when  $R = 0$ . Thus, the solutions are the level sets of  $E$ , which has a local minimum at  $\theta = 0, v = 1$ . Around a local minimum, the level sets are simple closed curves, so the fixed point is indeed a center.

$0 < R < 2\sqrt{2}$ : Here the fixed point (A) corresponds to a glider which dives with a constant velocity and not too steep an angle (the steepest angle is slightly less than  $\pi/4$  below horizontal). The fixed point is a spiral sink: nearby initial conditions spiral into it (B). Such a spirals corresponds to a glider for which the angle and velocity constantly oscillate, but these oscillations get smaller and smaller as time passes, limiting on the same angle and velocity as the fixed point. Solutions with initial conditions further away (D) from the fixed point do some number of loops before settling into a pattern of oscillations which limit on the diving solution. Between the solutions which do  $k$  loops and those which do  $k + 1$  loops lies a solution which stalls out after  $k$  loops (C) — at some time the solution gets to  $\theta = 2k\pi + \pi/2$  (i.e. going straight up) and velocity  $v = 0$ .

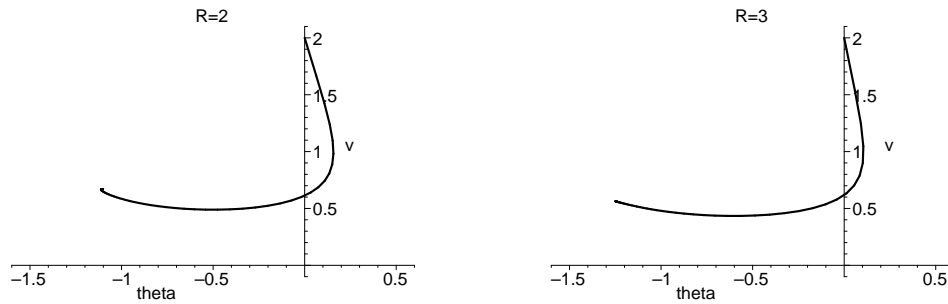


$R \geq 2\sqrt{2}$ : The fixed point is a sink corresponding to a steeply diving solution.<sup>6</sup> Causing the glider to loop becomes increasingly difficult as  $R$  increases. For example when  $R = 3$ , if the initial angle  $\theta$  is zero, an initial velocity of larger than 86.3 is required to get the glider to do one loop. If the initial angle differs from that of the fixed point, the glider fairly quickly turns towards that angle. No oscillation occurs; the angle can pass the limiting angle at most once.

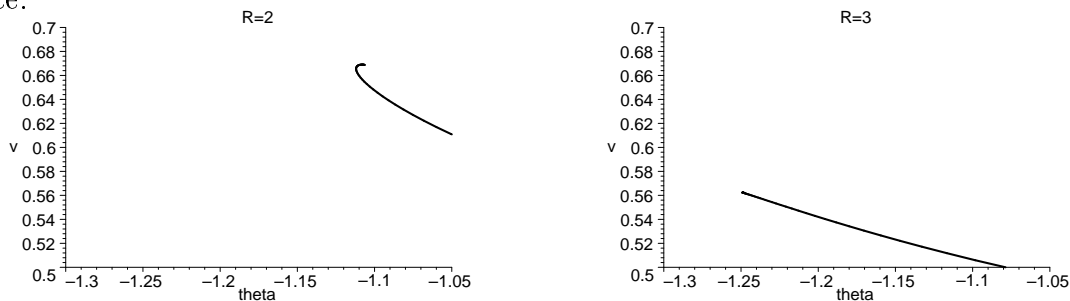


We have shown that for  $R < 2\sqrt{2}$ , the fixed point is a spiral sink: nearby solutions oscillate toward it. We should, however, point out that while the oscillations are always present mathematically, they can be very hard to discern. To emphasize this, we will compare a solution in the  $\theta, v$ -plane for  $R = 2$  (a spiral sink) and  $R = 3$  (a sink with two real eigenvalues).

<sup>6</sup>For  $R = 2\sqrt{2}$ , the linearization corresponds to a degenerate case with a double eigenvalue, rather than a regular sink. However, what we say here still applies.



The two plots look rather similar. However, if we look closely at the solutions near the fixed point (we can use zoom to do this without recomputing the picture), we see a significant difference.



For  $R = 2$ , there is a “hook” at the end, but for  $R = 3$ , the solution goes straight in to the fixed point. Further magnifications show the same pattern, as you may want to verify for yourself. The spiral for  $R = 2$  is always present, but very tight and hard to detect. Similarly, the oscillations in the  $x, y$ -plane when  $R$  is just slightly less than  $2\sqrt{2}$  are nearly indiscernable, but present.

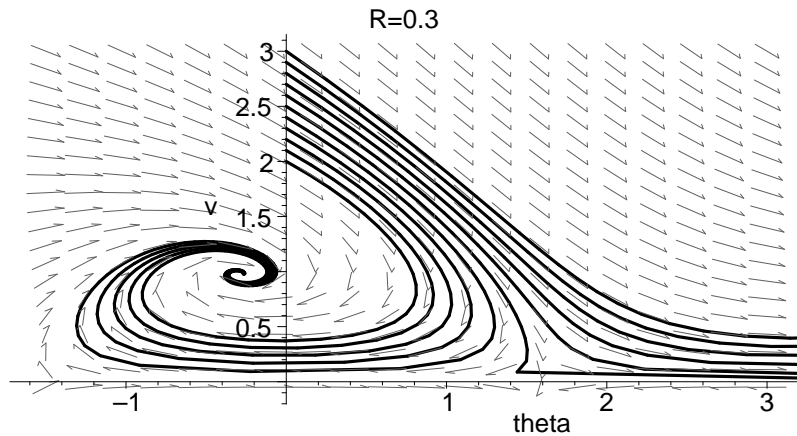
## 8 Dealing with the Singularity

The phugoid equations

$$\dot{\theta} = \frac{v^2 - \cos \theta}{v} \quad \dot{v} = -\sin \theta - Rv^2$$

run into trouble when  $v = 0$ , because  $\dot{\theta}$  becomes infinite as  $v \rightarrow 0$ . Since there are very few solutions for which  $v$  is ever zero, you might be tempted to ignore this issue, as we have done so far. But in fact, since we are solving the system numerically (rather than exactly), we run into problems even when  $v$  is merely small. So far, we have taken care to avoid such situations. But now let’s throw caution to the wind, and see what happens.

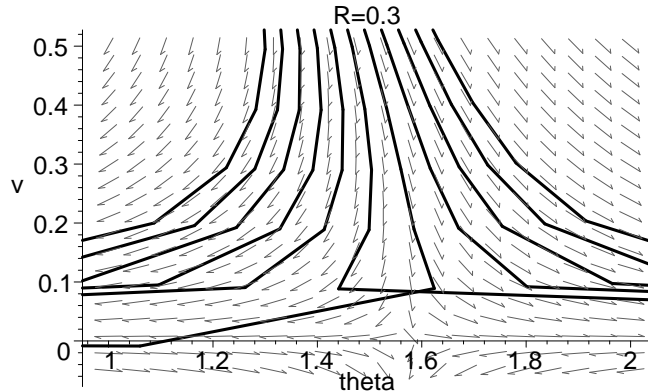
```
> R:=0.3:
DEplot([diff(theta(t),t)=(v(t)^2 - cos(theta(t)))/v(t),
diff(v(t),t) = -sin(theta(t)) - R*v(t)^2],
[theta,v], t=0..10,
[seq([theta(0)=0,v(0)=2+i/10],i=0..10)],
theta=-1.5..3, v=-0.05..3, title="R=0.3",linecolor=black,stepsize=0.1);
```



Something looks very wrong with the solution which starts at  $\theta = 0, v = 2.5$ . Notice how the solution makes a sharp turn, heading in a completely different direction from the vector field. Why is this happening?

Just to verify that it isn't a fluke, let's examine several solutions with initial conditions  $\theta(0) = 0$  and  $2.4 < v(0) < 2.6$ , looking in the problem area. We'll use `obsrange=false` to tell Maple we want it to continue computing the solutions, even when they go outside of the viewport.

```
> DEplot( [diff(theta(t),t) = (v(t)^2 - cos(theta(t)))/v(t),
           diff(v(t),t) = -sin(theta(t)) - R*v(t)^2 ],
  [theta,v], t=0..10,
  [seq([theta(0)=0,v(0)=2.4+.02*i],i=0..10)],
  theta=1..2, v=-0.05..0.5, obsrange=false,
  title="R=0.3",linecolor=black,stepsize=0.1);
```



Even worse! Not only does the solution for  $v(0) = 2.5$  go the wrong way, the one with  $v(0) = 2.52$  also turns the wrong way, crossing over another solution. Since our equations are autonomous, this is *never* supposed to happen.

We can remedy this by decreasing the stepsize significantly. Decreasing the stepsize by a factor of 10 (to 0.01) in the above example gives solutions which do the right thing. But a significant price in computation must be paid, and this doesn't really solve the problem. There are other initial conditions nearby which lead to the same trouble.

In order to decide how to fix the problem, we must first understand what is going wrong. The first problem we noticed happened for a solution with  $\theta(t_i) \approx 1.44$ ,  $v(t_i) \approx 0.09$ , using a stepsize  $h = 0.1$ . Since Maple is using a Runge-Kutta method as described in §4.2, let's calculate what is happening.

Recall that in Runge-Kutta, we evaluate the vector field at four points, and average them together to get the next point. We denoted these vectors as  $\vec{k}_\ell$ ,  $\vec{k}_c$ ,  $\vec{k}_m$ , and  $\vec{k}_r$ . Since we have to calculate the vector field at four points, we'll write a little function to do this for us, and use the fact that Maple lets us do arithmetic on lists as though they are actual vectors; that is, we can add them and multiply by scalars.<sup>7</sup>

```
> h:=0.1:
   VF:=p->[(p[2]^2 - cos(p[1]))/p[2], -sin(p[1])-R*p[2]^2]:

p:=[1.44, 0.9];
k[l]:=VF(p);
k[c]:=VF(p+h/2*k[l]);
k[m]:=VF(p+h/2*k[c]);
k[r]:=VF(p+h*k[m]);
step:=h/6*(k[l] + 2*k[c] + 2*k[m] + k[r]);

p := [1.44, 0.09]
k_l := [-1.359152319, -.9938883482]
k_c := [-4.858468860, -.9808007323]
k_m := [-8.872157510, -.9314790880]
k_r := [270.3559960, -.5250618133]
step := [4.025593182, -.08905849670]
```

The problem here is the vector  $\vec{k}_r$ , which has a huge  $\dot{\theta}$  component, and points in a different direction from the others. This happens because  $k_r$  is the value of the vector field calculated at  $\theta \approx 1.35$ ,  $v \approx -0.00315$ . Not only is the value of  $v$  very small here (giving a huge  $\dot{\theta}$  component), it is outside the range of allowed values, and the vector field points the wrong way. Similar problems will occur for any solution which comes too close to the  $v = 0$  axis with  $v(t)$  decreasing.

Now the question is how to avoid this? One not very satisfying answer is just to use a very small step size. But the problem with this idea is that we then must do a huge amount of computation, even when the vector field is behaving nicely.

A better solution is to use an *adaptive stepsize*: when the vector field is “nice”, we can use a large step, and when the vector field is being troublesome, we use a small step. Maple has such methods built in to it, such as the Runge-Kutta-Fehlberg method, which we can use by appending `method=rkf45` to the list of options to `DEplot`.

Sometimes, it is not always possible to change the numerical method. In the case of the phugoid, we can accomplish the same goal (in fact, gain something extra) by adjusting the

---

<sup>7</sup>This ability is not present in releases prior to Maple 6, making this computation a little more cumbersome.

length of the vectors that make up the vector field.

If we think of our solution curves as being parametric curves with tangent vectors which coincide with the vector field, then changing the lengths of the vectors (but not their direction) will not change the solution *curves*, only the parameterization. That is, we can rescale time so that the vectors no longer blow up as solutions approach the  $v = 0$  line.

Multiplying each vector by  $v(t)$  cancels out the problems in  $\dot{\theta}$ , and gives a system which has the same solution curves as long as  $v > 0$ . Thus, we get the new system

$$\dot{\theta} = v^2 - \cos \theta \quad \dot{v} = -v \sin \theta - Rv^3 \quad \dot{x} = v^2 \cos \theta \quad \dot{y} = v^2 \sin \theta$$

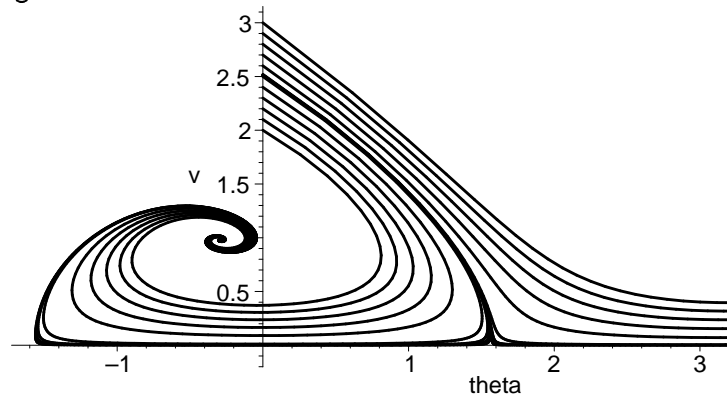
This new system of equations is defined and continuous for all values of  $v$  and  $\theta$ .

Furthermore, if we want to keep track of the original time variable  $t$ , we can add a new variable  $T$ . In the original, non-rescaled equations, the time could be described as  $T(t) = t$ , or equivalently,  $\dot{T} = 1$ . After rescaling the vector field by the factor of  $v$ , this gives  $\dot{T} = v$ .

```
> R:='R';
vphug:= [ diff(theta(t),t) = v(t)^2 - cos(theta(t)),
          diff(v(t), t)   = -v(t)*sin(theta(t)) - R*v(t)^3,
          diff(x(t), t)   = v(t)^2 *cos(theta(t)),
          diff(y(t), t)   = v(t)^2 *sin(theta(t)),
          diff(T(t), t)   = v(t) ];

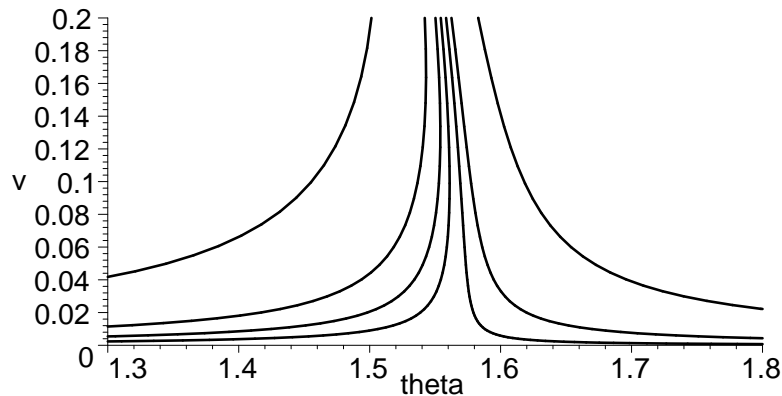
R:=0.3:

DEplot(vphug, [theta, v, x, y, T], t=0..20,
        [seq([theta(0)=0, v(0)=2+i/10, x(0)=0, y(0)=5, T(0)=0],i=0..10),
          [theta(0)=0, v(0)=2.51, x(0)=0, y(0)=5, T(0)=0],
          [theta(0)=0, v(0)=2.52, x(0)=0, y(0)=5, T(0)=0],
          seq([theta(0)=0, v(0)=2.51+i/1000, x(0)=0, y(0)=5, T(0)=0],i=2..5) ],
        theta=-1.5..3, v=-0.05..3, x=0..10, y=0..5, T=0..20,
        obsrange=false, linecolor=black, scene=[theta,v], stepsize=0.1);
```



```
> zoom(%, 1.3..1.8, 0..0.2);
```





As you can see, using the desingularized equations solves the problems for small values of  $v$ .

In addition, we can see additional structure not apparent in the original equations. There are new fixed points at  $v = 0$ ,  $\theta = \pm\pi/2$ . These fixed points are saddles, as we can see by looking at the Jacobian. In the desingularized equations, the Jacobian is

$$\begin{pmatrix} \sin \theta & 2v \\ -v \cos \theta & -\sin \theta - 3Rv^2 \end{pmatrix}$$

so the linearization at  $(\frac{\pi}{2}, 0)$  is  $\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$ , and at  $(-\frac{\pi}{2}, 0)$  we have  $\begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix}$ . These fixed points correspond to the solutions which stall out, and as we surmised earlier, they divide up the qualitative behaviors. The eigenvectors for both of them are parallel to the  $\theta$  and  $v$  axes, which says that solutions tending to a stall (or leaving one) do so by flying nearly straight up (resp. down) with an angle of  $\pi/2$  (resp.  $-\pi/2$ ). Solutions with very small velocity and an angle slightly less than  $\pi/2$  flip over rapidly to an angle slightly more than  $-\pi/2$  and then begin to increase their velocity.

While it is possible to deduce this behavior from the original equations, it is much more obvious once the singularity is filled in. Similar techniques are used in other branches of mathematics (blowing up a singularity in algebraic geometry, construction of a collision manifold in celestial mechanics) with great success.

## References

- [BDH] P. Blanchard, R. Devaney, G. R. Hall, *Differential Equations*. Brooks/Cole, 1998.
- [EP] C. H. Edwards, D. E. Penney, *Differential Equations & Linear Algebra*. Prentice-Hall, 2001.

- [HW] J. H. Hubbard, B. H. West, *Differential Equations: A Dynamical Systems Approach. Part 1: Ordinary Differential Equations*, Texts in Applied Mathematics **5**. Springer-Verlag, New York. 1991. *Part 2: Higher Dimensional Systems*, Texts in Applied Mathematics **18**. Springer-Verlag, New York. 1995.
- [Lan] F. W. Lanchester, *Aerial Flight Vol II: Aerodnetics*. Colchester & Co., London. 1908.

# fsqFsHn sGGousG

---

## Subsections

- [Introduction to Cryptography](#)
- [Simple Ciphers](#)
  - [Simple substitution](#)
  - [The Caesar cipher, and the ASCII encoding](#)
    - [Treating characters as numbers](#)
    - [Render unto Caesar](#)
- [Defining functions with proc; Local and global variables](#)
- [Caesar cipher redux](#)
- [Improved Caesar-like ciphers](#)
  - [The Vignère cipher](#)
  - [One-time pads](#)
  - [Multi-character alphabets](#)
- [Reading and Writing from a file](#)
- [Affine enciphering](#)
  - [When do affine encodings fail?](#)
  - [Implementing and using an affine encoding](#)
  - [Breaking an affine cipher](#)
- [Enciphering matrices](#)
  - [Treating text as vectors](#)
  - [Affine encoding with matrices](#)
  - [A Known-plaintext attack on an affine matrix cipher](#)
- [Modern cryptography](#)
  - [Secure cryptosystems](#)
  - [Message digests](#)
  - [Public Key cryptography](#)
- [Some Number Theory](#)
  - [The greatest common divisor and the Euclidean algorithm](#)
  - [The Chinese Remainder Theorem](#)
  - [Powers modulo  \$n\$](#)
  - [The Euler  \$\varphi\$ -function and Euler's Theorem](#)
- [The RSA Public key cryptosystem](#)
  - [Details of the RSA algorithm](#)
  - [Implementing RSA in Maple](#)
    - [Implementing the basics of RSA](#)

- [Making it Useful](#)
  - [RSA encoding a file](#)
- 

*Translated from LaTeX by Scott Sutherland*  
2002-08-29

# Chapter 4

## fsqF'sHn sGGousG

### 1 Introduction to Cryptography

If some information is meant to be kept private, the best means is to keep it well hidden. This is not, of course, always possible. One way around this is to hide the information or message in plain sight, that is, encode it in some way so that even if it is seen, it will be unreadable. The study of such encoding processes is called **Cryptography**<sup>1</sup>. The name “Cryptography” comes from the Greek words *kryptos* (κρυπτος), meaning “hidden” or secret, and *graphia* (γραφία), meaning writing.

Since ancient times, cryptography has been a part of military and governmental communications. More recently it has become part of nearly everyone’s life because of the Internet, electronic banking, and so on.

The usual jargon is as follows: the message you want to hide is called the **plaintext**, and the act of encoding it is called **encryption** or **enciphering**. The encoded plaintext is called the **crypttext** or the **ciphertext**, and the act of decoding it is called **decryption** or **deciphering** (or “cracking the code”, if the decoder wasn’t the intended reader). Usually, an encryption system (also called a cipher) has an auxiliary piece of information called the **key** needed for the encoding

---

<sup>1</sup>Cryptography is related to, but distinct from, **steganography**, which is the process of hiding the fact that a message exists at all. Using “invisible ink” is steganography, as is hiding messages within other messages. For example, looking at the second letter of each word in the message below (which was actually sent by a German spy during World War I [Kahn, p. 521]),

Apparently neutral’s protest is thoroughly discounted and ignored. Isman hard hit. Blockade issue affects pretext for embargo on byproducts, ejecting suets and vegetable oils.

one finds a rather different message, namely

Pershing sails from NY June I.

Steganography is often used to augment cryptography, and is also related to “digital watermarking”. We will not be considering steganography here.

and decoding process. Mathematically, we can represent the encryption process as

$$f_{key}(\mathcal{P}) = \mathcal{C},$$

where  $\mathcal{P}$  is the plaintext and  $\mathcal{C}$  is the ciphertext. To decrypt the message, one applies  $f^{-1}$  to the ciphertext. Note that to be able to decipher a message without ambiguity,  $f^{-1}$  must be a well-defined function (although  $f$  needn't be! It is perfectly all right to have the same message encode to two different ciphertexts, as long as we can get the original when deciphering.)

## 2 Simple Ciphers

### 2.1 Simple substitution

One of the most common (and very easy to crack) ciphers is substitution. One sometimes sees these in a newspaper somewhere near the crossword puzzle. This scheme is also the one used in Poe's story "The Gold Bug", and the Sherlock Holmes "Mystery of the Dancing Men". The idea is straightforward: choose a rearrangement of the letters of the alphabet, and replace each letter in the plaintext by its corresponding one.

Such substitutions are very simple to implement in Maple. We will use the `CharacterMap` function from the `StringTools` package<sup>2</sup> which does exactly what we want. First, we define two "alphabets"—one for our plain text, and one for our encrypted text. The second is a permutation of the first.

```
> with(StringTools):
  Alphabet := "abcdefghijklmnopqrstuvwxy":
  Cryptabet := "THEQUICKBROWNFXJMPSVLAZYG":
```

Now we define a pair of functions<sup>3</sup> to encrypt and decrypt our text.

```
> Scramble := msg -> CharacterMap(Alphabet, Cryptabet, msg):
  Unscramble := code -> CharacterMap(Cryptabet, Alphabet, code):
> Scramble("please do not read this");
```

"JWUTSU QX FXV PUTQ VKBS"

---

<sup>2</sup>The `StringTools` package was introduced in Maple 7, and the string data type was introduced in Maple V release 5. Everything we do in this chapter can easily be implemented without using `StringTools`—in fact, the original version of this chapter was written for Maple Vr4 using names rather than strings. Since Maple 6 is still in widespread use (Spring 2002), we'll try to include alternatives to using this package when practical.

<sup>3</sup>In Maple Vr5 and Maple 6, we could define `CharacterMap` as

```
CharacterMap := (old, new, txt) -> cat(seq(new[SearchText(txt[i], old)], i=1..length(txt)));
```

This will only work if all characters in the `txt` correspond to ones in `old`, so a space would have to be added to both `Alphabet` and `Cryptabet` to get this example to work.

```
> Unscramble(%);
```

“please do not read this”

Note that our message contains a spaces which are preserved in the encryption process, because the `CharacterMap` function only modifies those characters which are found in the first string. If a character isn’t found, it is left alone. Thus, if we try to scramble a message containing punctuation or upper-case letters, they aren’t modified by `Scramble`. Since our `Cryptabet` is all upper-case, a message which began with mixed case will be distorted.

```
> Scramble("He said I need direction; I head for the door.");
Unscramble(%);
```

“HU STBQ I FUUQ QBPUEVBXF; I KUTQ IXP VKU QXXP.”

“be said f need direction; f head for the door.”

It is traditional in pen-and-paper cryptography to use a 26-letter alphabet (that is, to remove spaces and punctuation and ignore any distinction between upper and lower case letters). If you wish to do this, it could be done<sup>4</sup> using `LowerCase` (which translates all letters to lower case), and a combination of `Select` and `IsLower` to remove anything that isn’t a letter.

```
> Select(IsLower,LowerCase("I don't mind suggestions. Don't give me any more."));
```

“idontmindsuggestionsdontgivemeanymore”

```
> Scramble(%);
Unscramble(%);
```

“BXFVCNBFQSLCCUSVBXFSQXFVCBAUNUTFDNXP”

“idontmindsuggestionsdontgivemeanymore”

In general we won’t adhere to this convention in this chapter. Rather, we will treat spaces and punctuation like any other character. This is not a common approach, but one we happen to like.

There are two big drawbacks to such a substitution cipher. One primary problem is that the key for enciphering and deciphering is the permutation: you need to remember an ordering of the entire alphabet.

---

<sup>4</sup>How to accomplish this without `StringTools` is left to the seriously motivated reader.

A second problem, shared with all such monoalphabetic substitutions (i.e., those that substitute one letter for another), is that they aren't that hard to break. If you have a long enough message, you can use the fact that certain letters occur much more frequently than others to help tease out the original text. In English text the letter "e" is the occurs most frequently, followed by "t", "a", "n", "i", and so on. Spaces occur about twice as often as the letter "e", assuming they are encrypted, rather than just omitted. Once you guess a few letters, common words like "the" and "and" fall into place, giving more clues.

We will examine some polyalphabetic substitution ciphers in §5.

## 2.2 The Caesar cipher, and the ASCII encoding

A variation on the arbitrary permutation discussed above is the "Caesar cipher", named after Julius Caesar, who supposedly invented it himself. This is also what was done by a "Captain Midnight Secret Decoder Ring", popular in the 1930s and 40s. Here we convert our alphabet to numeric equivalents (with, say A=0, B=1, and so on), add an offset to each numeric equivalent (legend has it that Caesar used an offset of 3), then re-encode the numbers as letters. We do the addition modulo the length of the alphabet; that is, when we shift a Y by 3, we get B, since  $24 + 3 = 27$ , and  $27 \equiv 1 \pmod{26}$ . However, we won't always restrict ourselves to a 26 letter alphabet—we may want punctuation, numbers, and to preserve upper- and lower-case distinctions.<sup>5</sup>

While we could easily implement the Caesar cipher without converting our text to numbers, doing so makes it easier to understand and sets us up for more complicated ciphers we will need later.

### Treating characters as numbers

Computers already have a built-in numeric representation of the alphabet—the ASCII encoding.<sup>6</sup> At their lowest level, digital computers operate only on numeric quantities—they just manipulate ones and zeros. The smallest piece of information is the state of a switch: it is off (0) or it is on (1). This is called a **bit**, which is short for **binary digit**. A group of 8 bits is called a **byte**, and is a single "character" of information.<sup>7</sup> A byte can have  $2^8 = 256$  different states:

00000000, 00000001, 00000010, 00000011, ..., 11111111

---

<sup>5</sup>For example, the title of this chapter (fsqFsHn sGGousG) is actually the result of using a Caesar cipher on its original title. A 53-character alphabet consisting of all the upper-case letters, a space, and all the lower-case letters was used, so the space in the middle may not correspond to a space in the original. The three Gs that occur should be a good hint to decoding it.

<sup>6</sup>In fact, while most computers these days use ASCII, it is not the only such arrangement. IBM mainframes, for example, primarily use another encoding called EBCDIC. But this makes no real difference to our discussion.

<sup>7</sup>The next larger unit in a computer is a **word**. However, different makes of computers have different sizes of words. Most computers today have 4 byte words (32 bits), but word sizes of 2 bytes (16 bits) and 8 bytes (64 bits) are not uncommon. 8-bit bytes are now essentially universal although in the early days of computing, there were computers which used other sizes, such as 6 bit bytes.



which are generally interpreted either as representing the integers 0 through 255, or the highest bit is often taken to represent the sign, in which case we can represent integers between  $-128$  and  $127$ . (We get an “extra” negative number because we don’t need  $-0$ , so the byte  $10000000$  can mean  $-128$  instead.) For notational convenience, bytes are often written in hexadecimal (that is, in base 16), grouping the bits into two blocks of 4 bits (sometimes called a *nybble*, half a byte), and using the characters A–F for the integers 10–15. Thus, the byte  $01011111$  would be written as  $5F$  in hexadecimal, since  $0101$  means  $2^2 + 2^0 = 5$ , and  $1111$  is  $2^3 + 2^2 + 2^1 + 2^0 = 15$ .

So how does a computer represent characters? There is an agreed upon standard encoding of characters, punctuation, and control codes into the first 127 integers. This is called *ASCII*,<sup>8</sup> an acronym for American Standard Code for Information Interchange. *Extended ASCII* encodes more characters into the range 128–255 (this is where you would find characters for å and ±, for example), although the characters found in this range vary a bit. Not all programs can deal successfully with characters in the extended range<sup>9</sup>

Maple has a built-in routine to allow us to convert characters to and from their ASCII equivalents. If it is called with a list of numbers, it returns a string of characters, and if called with a character string, out comes a list of numbers. For example,

```
> convert("How do I work this?",bytes);
```

```
[72, 111, 119, 32, 100, 111, 32, 73, 32, 119, 111, 114, 107, 32, 116, 104, 105, 115, 63]
```

```
> convert([72, 117, 104, 63],bytes);
```

“Huh?”

We can use this to make a list of the ASCII codes. Maple prints a  $\square$  when the ASCII code is not a printable character. We make the table for characters up through 127 because the meanings in extended ASCII differ depending on what character set you use. Feel free to make your own table of all 255 characters. (In the interests of making a nicely formatted table, we have used a some-

---

<sup>8</sup>ASCII is not the only such standard, although it is currently the most commonly used. Prior to the widespread use of personal computers in the 1980s, *EBCDIC* (Extended Binary Coding for Data Interchange and Communication) tended to be in more common use, because it was what IBM mainframes used. There are many other such assignments.

<sup>9</sup>In the late 1980s and early 1990s, there was a big push to make programs “8-bit clean”, that is, to be able to handle these “unusual” characters properly. There are a number of distinct assignments of characters to the upper 128 positions. One common one is *Latin-1* (ISO 8859-1), which contains the characters used in the majority of western European languages. Other assignments include characters for Scandinavian languages, Arabic, Cyrillic, Greek, Hebrew, Thai, and so on. Since the late 1990s, a multi-byte character encoding called *Unicode* (ISO 10646) has been in use; this universal character set allows computers to handle alphabets (such as Chinese) that need more than 255 characters. While use of Unicode is becoming more widespread, single byte character sets like ASCII will likely predominate for quite a long while.

what complicated command. The much simpler `seq([i,convert([i],bytes)],i=0..127)`; would have produced the same information, just not as neatly.)

```
> matrix([ [_ ,seq(i, i=0..9)],
           seq([10*j, seq(convert([i+10*j],bytes), i=0..9)], j=0..12)]);
```

-	0	1	2	3	4	5	6	7	8	9
0	"	"□"	"□"	"□"	"□"	"□"	"□"	"□"	"□"	"□"
10	"\\n"	"\\v"	"\\f"	"\\r"	"□"	"□"	"□"	"□"	"□"	"□"
20	"□"	"□"	"□"	"□"	"□"	"□"	"□"	"□"	"□"	"□"
30	"□"	"□"	" "	"! "	" \" "	"# "	"\$ "	"% "	"& "	" ' "
40	" (" "	" ) "	" * "	" + "	" , "	" - "	" . "	" / "	" 0 "	" 1 "
50	" 2 "	" 3 "	" 4 "	" 5 "	" 6 "	" 7 "	" 8 "	" 9 "	" . "	" , "
60	" < "	" = "	" > "	" ? "	" @ "	" A "	" B "	" C "	" D "	" E "
70	" F "	" G "	" H "	" I "	" J "	" K "	" L "	" M "	" N "	" O "
80	" P "	" Q "	" R "	" S "	" T "	" U "	" V "	" W "	" X "	" Y "
90	" Z "	" [ "	" \\ "	" ] "	" ^ "	" _ "	" ` "	" a "	" b "	" c "
100	" d "	" e "	" f "	" g "	" h "	" i "	" j "	" k "	" l "	" m "
110	" n "	" o "	" p "	" q "	" r "	" s "	" t "	" u "	" v "	" w "
120	" x "	" y "	" z "	" { "	"   "	" } "	" ~ "	" □ "	" □ "	" □ "

Note that some of the characters (such as ASCII code 10) print as a backslash and a letter, like `\\n`. These are codes for special control characters: `\\n` indicates a newline, `\\t` is a tab, and so on. Since Maple indicates a string by printing it inside double quotes, it indicates the double quote<sup>10</sup> character (ASCII 34) using `\"`, and the backslash character (ASCII 92) with a pair of backslashes. While it doesn't print out in Maple, ASCII 127 is the control sequence DEL (sometimes printed as `^?`). This sometimes is interpreted as "delete the previous character"; if your output device interprets it that way, characters can seem to mysteriously vanish.

Another special character is ASCII code 0: note that it prints as nothing, not even a blank space (ASCII 32 is the space character). This is the null character, and is typically used internally to indicate the end of a string. If a null occurs in the middle of a string, usually what comes after it won't print out unless you take special steps. To see this for yourself, try a command like `convert([73,116,39,115,0,103,111,110,101,33],bytes)`. You'll see that the characters following the 0 don't print out.

The `StringTools` package also contains `Char` which gives the character corresponding to a decimal number, and `Ord` which gives the number corresponding to a character. We'll generally use `convert`, although it makes little difference.

<sup>10</sup>The ASCII code does not have distinct codes for the open-quote (") and close-quote (") characters used in typesetting. Neither do most computer keyboards. However, the apostrophe (') and grave accent or backquote (`) characters have distinct ASCII codes (39 and 96, respectively).

## Render unto Caesar

Now we are nearly ready to write a version of the Caesar cipher. The basic idea is that for each character of the plaintext, we convert it a numerical representation, add a pre-specified offset modulo 127, then convert back to characters. Unlike the `Scramble` function we wrote in §2.1, we needn't declare an alphabet explicitly because we will use the ASCII codes. However, in addition to the Maple command `convert( ,bytes)`, we will use another new command, `map`.

The `map` command allows us to apply a function to every element of a vector or list in one call. For example, to compute the square root of each of the first 10 integers, we could do

```
> map(sqrt, [1,2,3,4,5,6,7,8,9,10]);
```

$$[1, \sqrt{2}, \sqrt{3}, 2, \sqrt{5}, \sqrt{6}, \sqrt{7}, 2\sqrt{2}, 3, \sqrt{10}]$$

If the function we wish to apply takes two variables, we can give the second one as a third argument to `map`. For example, the command

```
> map(modp, [1,2,3,4,5,6,7,8,9,10], 3);
```

$$[1, 2, 0, 1, 2, 0, 1, 2, 0, 1]$$

computes  $n \bmod 3$  as  $n$  ranges over the list  $[1, 2, \dots, 10]$ . (Recall that  $n \bmod p$  is the remainder after dividing  $n$  by  $p$ .)

We *finally* are ready. See if you can understand it before reading the explanation below. Remember to read from the middle outward. Note that we are defining the function we are mapping directly as the first argument.

```
> Caesar:= (plain, shift) ->
      convert(map(x -> (x+shift-1) mod 127 + 1,
                  convert(plain,bytes)),
            bytes):
```

Let's give it a try:

```
> Caesar("Et tu, Brute?", 3);
```

$$\text{"Hw\#wx/\#EuxwhB"}$$

If you don't see what is going on in `Caesar`, try to think of it as a composition of several functions, each of which you DO understand. That is, `Caesar` works as follows

- First, we convert the input string (called `plain`) into a list of numbers, using the call `convert(plain,bytes)`. Note that this is done on the third line, the *innermost* function.

- Next, we apply the function `x ->(x+shift-1) mod 127 + 1` to each of numbers in the resulting list using `map`. This call to `map` has two arguments: the definition of the shift, and what to apply it to (the innermost function `convert(plain,bytes)`).
- Finally, we convert the list back to characters again. This is the outermost function: it begins on the second line (immediately after the `->` and ends with the final parenthesis).

Viewing it as a composition, we have

$$\text{Caesar} := \text{convert}^{-1} \circ \text{map} \circ \text{convert}.$$

Note that we didn't really compute  $(x + \text{shift}) \bmod 127$ ; instead we subtracted one from the sum, computed  $\bmod 127$ , then added 1 again. Why? Consider what would happen if, for example, our character was "k" (ASCII code 107) and we shifted by 20. Then  $(107+20) \bmod 127$  gives 0, which is the code for the null byte, indicating the end of a string. This would be a problem. So instead, we add 1 afterward to ensure that our answer is in the range  $1, \dots, 127$ . And just to avoid confusing ourselves, we subtract the 1 first, so that shifting by 0 doesn't change anything.

Aside from the fact that this function is a bit confusing to read and offers essentially no cryptographic security, what else is wrong with it? Consider the following example:

```
> Caesar("Who put the benzedrine in Mrs. Murphy's Ovaltine?", 86);
```

```
“.?FvGLKvK?<v9<EQ<;I@E<v@Ev$IJ□v$LIG?P}Jv&M8CK@E< □”
```

In addition to the fact that the result looks funny, if we were to try to decipher this from printed output, we would have a little (but not much) trouble, because the characters “.” and “?”, after encoding, print out as □.<sup>11</sup> This happens because the numeric values of these characters, after shifting by 86, are 5 and 22, respectively. Neither of those corresponds to a printing character, so we can't tell them apart. Of course, if we had no intention of dealing with our encoded message in printed form, this would be no problem. The actual character values are different, they just print out the same.

How can we get around this? We can restrict ourselves to printing characters using a variation of the idea in `Scramble` of §2.1, that is, giving an explicit `Alphabet` that we will deal in.

Before doing this, however, it will be helpful to discuss the `proc` command, which allows us to construct functions that consist of more than one Maple command.

---

<sup>11</sup>Of course, it is easy to make examples where significant parts of the message all come out as □. One such is the same message with a shift of 50.

### 3 Defining functions with proc; Local and global variables

We already know how to define functions in Maple. However, all of the functions we have defined so far have consisted of a single Maple expression (no matter how complicated it might be). For example,

```
> sqr := x -> sqrt(x);
```

```
sqr := sqrt
```

```
> sqr(25);
```

5

However, Maple has a more versatile way of defining functions using the `proc` command. The function `sqr2` below behaves exactly the same as `sqr` defined above. (When using `proc`, the function begins with the keyword `proc`, followed by any parameters in parenthesis, then one or more Maple statements, and finally the keyword `end`. The result of the function is the result of the last Maple command executed.

```
> sqr2 := proc(x)
    sqrt(x);
end;
```

```
sqr2 := proc(x) sqrt(x) end
```

```
> sqr2(25);
```

5

```
> sqr2(-25);
```

5 I

Unlike the functions defined by `->`, however, we can use several statements when using `proc`:

```
> sqr3 := proc(x)
    if (x >= 0) then
        sqrt(x);
    else
        print("The square root of ",x,"is not a real number");
    fi;
```

```
end;
```

```
sqr3 := proc(x)
  if 0 ≤ x then sqrt(x) else print("The square root of ", x, " is not a real number") fi
end
```

```
> sqr3(-25);
```

*The square root of , -25, is not a real number*

```
> sqr3(25);
```

5

Here is another variation on the same idea: we compute a variable  $y = x^2$ . If  $y > 25$ , the result is 25, otherwise it is  $y$ .<sup>12</sup>

```
> mesa := proc(x)
  y := x^2;
  if ( y>25) then
    25;
  else
    y;
  fi;
end;
```

Warning, 'y' is implicitly declared local

```
mesa := proc(x) local y; y := x^2; if 25 < y then 25 else y fi end
```

What does that warning mean? It means that Maple has assumed that the variable  $y$  used in `mesa` is local to that procedure. That is, it is different from a variable  $y$  that may be used outside the context of the procedure. Let's check that. First, we set  $y$  to 35, invoke `mesa(19)` (which sets *its* copy of  $y$  to 361), then check what the value of  $y$  is.

```
> y:=35;
```

$y := 35$

---

<sup>12</sup>Note that this could also be done using `piecewise`.

```

> mesa(19);

25

> y;

35

```

It is still 35. Now, let's add a line saying that `y` has a *global* definition— the `y` within the procedure is the same as the `y` outside.

```

> mesa2:= proc(x)
  global y;

  y := x^2;
  if ( y>25) then
    25;
  else
    y;
  fi;
end;

  mesa2 := proc(x) global y; y := x^2 ; if 25 < y then 25 else y fi end

> y;

35

> mesa2(19);

25

> y;

361

```

This time, of course, `y` was indeed changed by invoking `mesa2`.

The opposite of `global` is, not surprisingly, `local`. In order to stop Maple from giving warnings about variables being implicitly declared local, we can add the statement `local y;` just after the `proc` statement. Using such statements is called **declaring the scope** of the variables. This is a good habit to get into, because it lessens the chance of accidentally using a

global variable or of misspelling the name of a variable. Some computer languages require that you declare all variables you use.

Another good habit to get into is to indicate what type of arguments the function will accept. This is optional in Maple, and not always desirable (you may not always know what they will be). This is done by specifying the type after two colons in the argument list. For example,

```
> ith:=proc(l::list, i::posint)
    return(l[i]);
end:
```

is a function that insists its first argument be a list, and the second must be a positive integer. If you call it with something else, you will get an error message.

```
> ith("Henry",8);
```

```
Error, invalid input: ith expects its 1st argument, l, to be of type list, but
received Henry
```

## 4 Caesar cipher redux

Now that we know about `proc`, let's implement the Caesar cipher again, this time with an alphabet of our choosing.

First, we define a global variable `Alphabet` that lists all the characters we intend to encode. We would like to be able to distinguish all characters which print out, as well as blanks (that is, those which you have keys for on your keyboard). To avoid typing in the full list of characters we want<sup>13</sup>, we generate a list of all the integers less than 256 and convert it to a character string with `convert(,bytes)`. Then we use `Select` and `IsPrintable` to pick out the usual printable characters. Finally, since we would like to have a space and a newline as characters in our alphabet, we add them to the resulting string with `cat`.

```
> Alphabet := cat("\n\t", Select(IsPrintable,convert([seq(i,i=1..255)],bytes)));
```

```
“\n\t !”#%&'()*+,-./0123456789;<=>?@ABCDEFGHIJKLMNPOQRSTUVWXYZ\
YZ[\]^_‘abcdefghijklmnopqrstuvwxyz{|}~”
```

We will use this particular alphabet regularly in the next few sections.

Next, we define two functions which perform the analogue of `convert(,bytes)`, but with our special alphabet. `ListToString` will take a list of integers and replace it with the character from that position in `Alphabet`, and `StringToList` does the reverse process. For each character in the input string, `StringToList` uses `SearchText` to find its position in `Alphabet`. We need

---

<sup>13</sup>For versions prior to Maple 7, just type in the desired alphabet directly



to subtract one from this position because we would like our characters to be numbered from 0 rather than 1. Thus, in the `Alphabet` above, `\t` is character 0, and the numeric code for `~` is 96. `Alphabet` is declared as a global variable.

```
> StringToList := proc(text::string)
    local i;
    global Alphabet;
    [seq(SearchText(text[i],Alphabet)-1, i=1..length(text))];
end:

> ListToString := proc(numlist::list(nonnegint))
    local i;
    global Alphabet;
    cat(seq(Alphabet[numlist[i]+1], i=1..nops(numlist)));
end:
```

Now the Caesar cipher is written in terms of these functions. Note that since 0 corresponds to a printable character, we don't have to add and subtract 1 as before.

```
> Caesar2:= proc(plaintext::string, shift::integer)
    local textnum,codenum,i,p;
    global Alphabet;

    p      := length(Alphabet);
    textnum := StringToList(plaintext);
    codenum := [seq( modp(textnum[i]+shift, p), i=1..length(plaintext)) ];
    ListToString(codenum);
end:
```

It works as follows.

- We are given a `plaintext` message, and a shift amount. The temporary variables that we will use during the conversion process are declared as local, and the `Alphabet` is declared as global.
- The number of characters in the alphabet is saved as `p`.
- The plaintext is converted to a list of numbers called `textnum`.
- Each number in `textnum` is shifted, and the result is taken mod `p`. Note that since `StringToList` yields result in  $\{0, 1, \dots, p - 1\}$ , we don't have to monkey around to get the result into the proper range like we did in the first `Caesar` (§2.2, p. 4:7). Also note that we use `modp(a,p)` to compute  $a \bmod p$ , instead of writing `a mod p` as before. Either way works fine.
- Finally, the result is converted back to characters with `ListToString`.

Let's try it out.

```
> Caesar2("Veni, Vidi, Vici",3);
```

```
“Yhql/#Ylgl/#Ylfl”
```

To decode, we can just use the negative of the shift amount.

```
> Caesar2(%, -3);
```

```
“Veni, Vidi, Vici”
```

```
> text:="I have heard the mermaids singing, each to each.
I do not think that they will sing to me.";
```

```
text := “I have heard the mermaids singing, each to each.\nI do not think \
that they will sing to me.”
```

In this sample, we have used a text which contains a newline character. Maple represents this as `\n` inside a text string. Since we will treat the newline as a regular character, it too will encode to something. If you look carefully at the output below, you can probably pick it out.

```
> Caesar2(text,64);
```

```
“(’G@UD’GD@QC’SGD’LDQL@HCR’RHMFMFI’D@BG’SN’D@BGn^\  
(’CN’MNS’SGHMJ’SG@S’SGDX’VHKK’RHMFM’SN’LDn”
```

```
> Caesar2(%, -64);
```

```
“I have heard the mermaids singing, each to each. \nI do not think that they \  
will sing to me.”
```

If we would like to see the decrypted text with the newlines expanded, we need to use `printf`.

```
> printf(%)
```

```
I have heard the mermaids singing, each to each.
I do not think that they will sing to me.
```

One thing that we should check is what happens if our message should contain a character that is not in our alphabet. Rather than changing our text, we’ll dramatically shorten the alphabet.

```
> Alphabet:="abcdefghijklmnop_";
```

```
Alphabet := “abcdefghijklmnop_”
```

Now let us go then, you and I, and re-encode the test message using the new alphabet, but with a shift of zero. Ordinarily, a shift by 0 would not change the text at all. But with the shortened alphabet, something happens to the characters that aren’t mentioned.

```
> Caesar2(text,0);
“_h_a_e_h_e_a_d__h_e_m_e_m_a_i_d__i_n_g_i_n_g__e_a_c_h__o__e_a_c_h____d_o_n_o__h_i_n_k__h_a__h_e__i_l_l__i_n_g__o_m_e_”
```

Can you explain why this happens?

## 5 Improved Caesar-like ciphers

Certainly the Caesar cipher offers no cryptographic security at all: if you know the alphabet the message was encoded in, you need only guess one character to crack the code. Even if you don't know the alphabet, guessing the correspondence is not very hard with a little patience.

In this section, we will discuss a few approaches to improving the security, while retaining the basic idea of character shifting.

### 5.1 The Vignère cipher

One way to make a Caesar cipher a bit harder to break is to use different shifts at different positions in the message. For example, we could shift the first character by 25, the second by 14, the third by 17, and the fourth by 10. Then we repeat the pattern, shifting the fifth character by 25, the sixth by 14, and so on, until we run out of characters in the plaintext. Such a scheme is called a Vignère cipher<sup>14</sup>, which was first used around 1600, and was popularly believed to be unbreakable.<sup>15</sup> This cipher is called a polyalphabetic substitution cipher, because several different substitutions are made depending on the position of the character within the text.

In our first example, the key consists of the four shifts [25, 14, 17, 10], which are the numerical equivalents of the string “ZORK” in a 26-letter alphabet consisting of the letters A–Z. It is common practice to think of our key as plaintext letters, rather than their numerical equivalents, but either will do. We can encode the string “CRYPTOGRAPH” as

<i>C</i>	<i>R</i>	<i>Y</i>	<i>P</i>	<i>T</i>	<i>O</i>	<i>G</i>	<i>R</i>	<i>A</i>	<i>P</i>	<i>H</i>
+25	+14	+17	+10	+25	+14	+17	+10	+25	+14	+17
<i>B</i>	<i>F</i>	<i>P</i>	<i>Z</i>	<i>S</i>	<i>C</i>	<i>X</i>	<i>B</i>	<i>Z</i>	<i>D</i>	<i>Y</i>

Note that in the above, the letter “R” in the plaintext encodes to both “F” and “B” in the crypttext, depending on its position. Similarly, the two “Z”s in the crypttext come from different plain characters.

---

<sup>14</sup>This cipher takes its name after Blaise de Vignère, although it is actually a corruption of the one he introduced in 1585. Vignère's original cipher changed the shift amount each letter based on the result of the last encoding, and never repeated. This scheme is *much* harder to break. However, one reason for its lack of popularity was probably due to the fact that a single error renders the rest of the message undecipherable. More details can be found in [Kahn].

<sup>15</sup>In fact, as late as 1917, this cipher was described as “impossible of translation” in a respected journal (Scientific American), even though the means to break it had been well known among cryptographers for at least 50 years.

Now we implement this in Maple. This is very similar to the Caesar cipher, just with the extra complication of multiple shifts, and letting our key be a string.

First, we set our Alphabet to the usual one. We also use the conversion functions from §4.

```
> Alphabet := cat("\n\t", Select(IsPrintable,convert([seq(i,i=1..255)],bytes)));
> Vignere:= proc(plaintext::string, key::string)
    local textnum,codenum,i,p,offsets,keylen;
    global Alphabet;

    p      := length(Alphabet);
    offsets := StringToList(key);
    keylen  := length(key);
    textnum := StringToList(plaintext);
    codenum := [seq( modp(textnum[i] + offsets[modp(i-1,keylen)+1], p),
                    i=1..length(plaintext)) ];
    ListToString(codenum);
end;
```

To try it out, we'll use the same text as in the previous section. Notice how much harder it is to pick out the word boundaries in the resulting ciphertext.

```
> coded:=Vignere(text,"Prufrock");

coded := "{t^HiUeT6ThKtdLQR[Y'QMPDtiPaWMZ8\twLTSLmEbwLTSL \
{Rr?hW_eZ@gw[[YRWRg^HgqXT6lw^\\\PmD\\dNtdSm>X$"
```

We can make the decoding function from the original (let's call it `unVignere`) by changing exactly one  $+$  sign to a  $-$ .<sup>16</sup> We omit the change here so perhaps you will figure it out for yourself. But we will test it, to show you that it does work.

```
> printf(unVignere(coded,"Prufrock"));

I have heard the mermaids singing, each to each.
I do not think that they will sing to me.
```

Even though this scheme looks quite daunting, it is not so very hard to crack if you use a computer or have a very large supply of perseverance. If we know that the key is of a certain length, say 4, and our plaintext is sufficiently long, then we can perform frequency analysis on every fourth letter. Even if we don't know the key length, it is not too hard to write a computer program to try all the lengths less than, say, 10, and pick the one that looks best.

## 5.2 One-time pads

Note that the longer the key is in the Vignère cipher, the harder it is to break. If the key were as long as the text, then it might seem at first that analyzing the frequency of letters in the encrypted text would be of no help, since each letter would be shifted by a different

---

<sup>16</sup>Note that we could also use the `Vignere` routine, but with the inverse of the key. For example, in the preceding example, the inverse of "Prufrock" is "M+(7+.:2": the numeric code of P plus the numeric code of M is 97 (the length of the alphabet), similarly for r and +, u and (, and so on.

amount. This is *almost* true. If the key is an passage of text in English, then the shifts will occur with a predictable frequency. Of course, the problem gets very difficult, but cryptanalysts are persistent people.

But what if there were no predictability within the key, having the shifts come at random? The result (a Vignère cipher with an infinitely long, completely random key) is a cryptosystem that *cannot be broken*. Since the shifts are random, if you manage to decipher part of the message, this gives you no clue about the rest. Furthermore, any plaintext of a given length can encrypt to any ciphertext (with a different key, of course). For example, the ciphertext “=5nwhn KDNO?uWTBC-XA” might have come from the phrase “Let’s have dinner.” (with the key “pOyntmbbXYtrjSTGe1”), or it might be the encryption of “Attack at midnight” with the key “{@y4#!Jbz>&moSYEoL”. Since any message could encrypt to any other, there is no way to break such a code unless you know the key.

But that is the problem: the key is infinitely long. Infinitely long, truly random sequences of numbers tend to be somewhat unwieldy. And to decode the message, you must know what random sequence the message was encoded with.

Such a system is called a **one-time pad**, and was used regularly by spies and military personnel. Agents were furnished with codebooks containing pages and pages of random characters. Then the key to the encryption is given by the page on which to begin. It is, of course, important that each page be used only once (hence the name “one-time pad”), because otherwise if a codebreaker were able to intercept a message and (via some other covert means) its corresponding translation, that could be used to decipher messages encoded with the same page. This sort of setup makes sense if an agent in the field is communicating with central command (but not with each other). Each agent could be given his own codebook (so that if he is captured, the whole system is not compromised), and he uses one page per message. Central command has on file the books for each agent.

A variation on this theme is the **Augustus cipher**,<sup>17</sup> where instead of a random sequence of shifts, a phrase or passage from a text which is as long as the plaintext is used. The trouble with this is that, because of the regularities in the key, a statistical analysis of the crypttext allows one to break the cipher.

Another issue is that to be truly unbreakable, the random sequence must be truly random, with no correlation among the characters. This is harder than it sounds—real randomness is hard to come by. If the random sequence has some predictability, the resulting stream can be attacked. A number of attacks on cryptosystems have been made not by breaking the encryption scheme directly, but because the underlying random-number generator was predictable.

We can easily modify our Vignere program to be a one-time pad system, using Maple’s random number generator to make our one-time pad.<sup>18</sup> You might think that generating a

---

<sup>17</sup>Sometimes a Caesar cipher with a shift of +1 is also called an “Augustus Cipher”, even though these are very different ciphers.

<sup>18</sup>Technically speaking, this is not a one-time pad, but a one-time stream. The distinction is subtle, and we will ignore it here.

random sequence of numbers would be inherently unreproducible, so the message would be indecipherable unless we record the stream of random numbers. However, computers can not usually generate a truly random sequence. Instead, they generate a **pseudo-random** sequence  $s_1, s_2, s_3, \dots$  where the pattern of the numbers  $s_i$  is supposed to be unpredictable, no matter how many of the values of the sequence you already know. However, to start the sequence off, the pseudo-random number generator requires a **seed**—whenever it is started with the same seed, the same sequence results. We can use this to our advantage, taking the seed to be our key.<sup>19</sup> Note that in order to decode the message by knowing the key (the seed), the recipient must use the *same* pseudo-random number generator.

Maple's random number generator gives different results when called in different ways. If called as `rand()`, it gives a pseudo-random, non-negative 12 digit integer. When called as `rand(p)`, it gives a procedure that can be called later to generate a random integer between 0 and  $p$ ; it is this second version that we will use. The seed for Maple's random number generators is the global variable `_seed`.<sup>20</sup>

```
> OneTimePad := proc(plaintext::string, keynum::posint)
    local textnum, codenum, i, p, randnum;
    global Alphabet, _seed;

    p      := length(Alphabet);
    randnum := rand(p);
    _seed  := keynum;
    textnum := StringToList(plaintext);
    codenum := [seq((textnum[i] + randnum()) mod p, i=1..length(plaintext))];
    ListToString(codenum);
end;
```

In this implementation, it is assumed that the key is a positive integer (it can be as large as you like). It would be easy to change it to use a string of characters, however, by converting the string to a number first. One way to do that is discussed in the next section.

In most descriptions of one-time pad systems, one takes the exclusive-or (XOR) of the random number and the integer of the plaintext, rather than adding them as we have done. This does not significantly change the algorithm if the random sequence is truly random, since adding a random number is the same as XOR'ing a different random number. However, the XOR approach has the advantage that the enciphering transformation is its own inverse. That is, if we produce the crypttext using `crypt:=OneTimePadXOR(plain,key)`, then `OneTimePadXOR(crypt,key)` will give the decryption with the same key. This is not the case

---

<sup>19</sup>Pseudo-random number generators appropriate for cryptography are rare. Most implementations (including Maple's) are good enough for everyday use, but not enough to be cryptographically secure. By analyzing the output of a typical random number generator, a good cryptanalyst can usually determine the pattern. For example, Maple's `rand` function (and that of most computer languages, such as C, Fortran, and Java) gives the result of an affine sequence of numbers, reduced to some modular base. That is,  $x_i = ax_{i-1} + b$  and  $s_i = x_i \bmod n$  for some fixed choices of  $a$ ,  $b$ , and  $n$ . In this setting, the seed is  $x_0$ . We shall ignore the problem that this sequence is guessable, but if you want real security, you cannot.

<sup>20</sup>We can choose a "random" seed (based on the computer's clock) using the function `randomize()`.

for the version given above; to make a decryption procedure, we would need to modify the above by changing the  $+$  to a  $-$ .

### 5.3 Multi-character alphabets

We can also improve security a bit by treating larger chunks of text as the characters of our message. For example, if we start with the usual 26-letter alphabet A–Z, we can turn it into a 676-letter alphabet by treating pairs of letters as a unit (such pairs are called **digraphs**), or a  $26^3$ -letter alphabet by using trigraphs, or triples of letters. This makes frequency analysis much harder, and is quite easy to combine with the other cryptosystems already discussed. We will use 99-graphs on a 256-letter alphabet (the ASCII code) when we implement the RSA cryptosystem in §11.2. While frequency analysis is still possible (charts of digraph frequencies are readily available, trigraphs less so), the analysis is much more complex.

To convert the digraph “HI” to an integer (using a length  $26^2$  alphabet of digraphs), one simple way is to just treat it as a base-26 number. That is, “HI” becomes  $7 \times 26 + 8 = 190$ , assuming the correspondence of H=7, I=8. To convert back, we look at the quotient and remainder when dividing by 26. For example,  $300 = 26 \times 11 + 14$ , yielding “LO”.

Either we can do this arithmetic ourselves directly, or we can use the `convert(,base)` command. This command takes a list of numbers in one base and converts it to another. One slightly confusing fact is that in this conversion, the least significant figure comes first in the list,<sup>21</sup> instead of the usual method of writing numbers with the most significant digit first.

For example,  $128 = 1 \times 10^2 + 2 \times 10^1 + 8 \times 10^0$  would be written as `[8,2,1]`. To convert this to base 16, we would note that  $128 = 8 \times 16^1 + 0 \times 16^0$ , so in base 16, it is written as 80.

Doing this calculation in Maple, we have

```
> convert([8,2,1], base, 10, 16);
```

[0, 8]

Below is one way to implement the conversion of text to numeric values for  $k$ -graphs. We assume our usual functions `StringToList` and `ListToString` are defined (see §4), as well as the global `Alphabet`. The routine below converts `text` into a list of integers, treating each block of  $k$  letters as a unit. A block of  $k$  characters  $c_1c_2c_3 \dots c_k$  is assigned the numeric value  $\sum_{i=0}^{k-1} x_i p^k$ , where  $x_i$  is the numeric equivalent of  $c_{i+1}$  assigned by `StringToList`.

---

<sup>21</sup>Such a representation is called little-endian, as opposed to a big-endian one. Some computers represent numbers internally in big-endian format (Sun SPARC, Power Macintosh), and others use a little-endian representation (those using Intel processors, which is what MS-Windows and most versions of Linux run on). This name comes from Gulliver’s Travels, where in the land of *Blefuscu* there is war between the big-endians (who eat their eggs big end first), and the little-endians (who eat eggs starting from the little end).

```

> StringToKgraph := proc(text::string, k::posint)
    local p;
    global Alphabet;
    p:= length(Alphabet);
    convert(StringToList(text), base, p, p^k);
end:

> KgraphToString := proc(numlist::list(nonnegint), k::posint)
    local p;
    global Alphabet;
    p:=length(Alphabet);
    ListToString( convert(numlist, base, p^k, p));
end:

```

In the examples below, we are using our usual 97-character alphabet. Of course, this will work on any alphabet, although the specific numbers will differ.

```

> StringToKgraph("two by two",2);

```

[8719, 275, 8895, 8344, 7946]

In our alphabet, “t” is character number 86 and “w” is number 89, so the digraph “tw” encodes as  $86+89 \times 97 = 8719$  (remember we are using a little-endian representation). Similarly, “o ” gives  $2 \times 97 + 81$ , and so on. Notice that the two occurrences of “two” give different numbers, because one begins on an even character and the other starts on an odd one. The first corresponds to the digraphs “tw” and “o ”, while the second is “ t” and “wo”.

We can also encode with 4-graphs, if we like.

```

> StringToKgraph("two by two",4);

```

One advantage (for us) of Maple’s use of the little-endian order is that we needn’t worry whether the length of our text is divisible by  $k$ . In the above example, the last 4-graph is “wo”, which is encoded as though it had two more characters on the end with numeric code of 0. The disadvantage of this is that if our text ends with the character with code 0 (in our standard 97-character alphabet, this is a newline), that character will be lost.

Another way to treat multiple characters together is to think of them as vectors. For example, the digraph “by” might correspond to the vector [68, 91]. We will treat this approach in §8.

## 6 Reading and Writing from a file

While none of the methods we have covered yet are secure enough to foil a good cryptanalyst, the OneTimePad procedure of §5.2 is good enough for household use (unless you live with



cryptanalysts!) Should we want to really use it, however, we would find it extremely tedious to type in our message each time, and then copy the encrypted text from the worksheet.

Fortunately, Maple is able to read and write data from a file, using `readline` and `writeline`. These read (or write) a single line of data from a specified file.<sup>22</sup> When reading from a file, `readline` returns either a line of data, or 0 when the end of the file is reached.

Although it is not strictly necessary, it is a good idea to explicitly `open` the file before dealing with it, and `close` it when you are done. In addition to being good practice (like explicitly declaring variables), it also allows you to specify certain options, like whether you intend to read or write to the file. In fact, I recommend using `fopen`, which also allows you to specify whether the file should be treated as text or binary.<sup>23</sup> Below is a small example of a procedure which encrypts a file one line at a time using a procedure `EncryptLine`. We have never given such a procedure, but it should be clear how to modify it to use your favorite encryption routine instead.

```
> Cryptfile := proc(plainfilename,cryptfilename,key)
    local line, cline, file, cryptf;

    file :=fopen(plainfilename,READ,TEXT);
    cryptf:=fopen(cryptfilename,WRITE,TEXT);

    line :=readline(file);
    while (line <> 0) do
        writeline(cryptf,EncryptLine(line,key));
        line :=readline(file);
    od;
    close(file);
    close(cryptf);
end;
```

## 7 Affine enciphering

We have just examined a trivial encryption scheme, the Caesar cipher, which can be described as applying  $x \mapsto x + b \pmod{p}$  to each of the characters  $x$  in a message. As we have seen, this provides no security at all (a message encoded in this way can be cracked in a short while by a third grader), but minor modifications such as those in §5 can greatly improve the security.

We now return to trivialities, but make it one step harder. Instead of encoding each character by just adding a constant, we can multiply by a constant as well. That is, we can use  $x \mapsto$

---

<sup>22</sup>If you are treating a newline as a character, you will need to deal with the file on a byte-by-byte basis, using `readbytes` and `writebytes`

<sup>23</sup>On several operating systems, including MS-Windows and MacOS (but not Unix or Linux), files may be treated as containing text or binary data. When dealing with a text file, some control characters (such as the end-of-line marker) are translated to a format more suitable for use. When opening the file, we can tell Maple whether it should be dealt with as text or binary. If you neglect to specify which, you may find your file has one very long line if you open it in a text editor. Under Unix, text and binary files are the same, so no translation is necessary, but there is no harm in specifying the file type.

$ax + b \pmod{p}$  to encode each character. This not only helps a tiny bit (to crack it, you need to find out *two* characters instead of just one, and the encoded message looks a little more “mixed up”, since letters adjacent in the alphabet no longer encode to adjacent letters), it helps us set up a little more conceptual machinery.

## 7.1 When do affine encodings fail?

You might be wondering if this will work at all. Certainly when we were just shifting the letters by a fixed amount, there was no chance that the enciphering transformation  $f$  could fail to be invertible. Does the same hold true if we multiply the numeric codes for the letters? That is, do different letters always encode to different letters?

The answer is, unfortunately, not always. For example, suppose we use a 27 character alphabet (A-Z and blank). Then if we encode using the transformation  $x \mapsto 6x + 2 \pmod{27}$ , both “A”(0) and “J”(9) encode to the character “C”(2). What went wrong? Recall that when we work modulo 27, any numbers that differ by a multiple of 27 are equivalent. Since 27 has nontrivial factors, several characters will collide whenever the multiplicand (6) shares a common divisor with 27. This will never happen if  $\gcd(a, p) = 1$ .

Let’s make that precise:

**Proposition 7.1** *Let  $a, b, x,$  and  $y$  be integers in the set  $\{0, 1, \dots, p - 1\}$ , with  $\gcd(a, p) = 1$ . Then*

$$ax + b \equiv ay + b \pmod{p} \iff x \equiv y \pmod{p}$$

**Proof.** Suppose  $ax + b \equiv ay + b \pmod{p}$ . We want to show that  $x \equiv y \pmod{p}$ . Certainly we have  $ax \equiv ay \pmod{p}$ , so there must be integers  $k$  and  $m$  so that  $kp + ax = mp + ay$ . But we can then rewrite this as

$$a(x - y) = (k - m)p$$

that is,  $a(x - y)$  is a multiple of  $p$ . If  $a$  and  $p$  share no common divisors, then  $x - y$  must be a multiple of  $p$ , so  $x \equiv y \pmod{p}$ .

If  $x \equiv y \pmod{p}$ , then certainly  $ax + b \equiv ay + b \pmod{p}$ . □

Thus, by the above, we can see that the affine encoding scheme will be just fine if the multiplicand  $a$  is relatively prime to the length of our alphabet ( $p$ ). Henceforth, we will always take an alphabet whose length is a prime number, so we need not worry about this issue. An alphabet of length 97 works quite well, because 97 is prime and large enough so that we can include all the usual printable ASCII characters (see §2.2, especially the table on page 4:6). If you plan to use the full ASCII or extended ASCII character set (including the null character), then you will need to make sure you choose  $a$  to be an odd number (since the length of the alphabet will be either 128 or 256, both of which are powers of 2).

## 7.2 Implementing and using an affine encoding

Implementing this scheme is a trivial modification to the `Caesar2` procedure we wrote in §4. We only need add another parameter `a`, and change the encoding calculation. We also renamed `shift` as `b` to correspond to our discussion here.

In light of the Prop. 7.1, we will also add some error checking to ensure that `a` and the length of the alphabet are relatively prime. Recall that if they are not, the encryption will produce a message that cannot be decrypted. We use the `error` command<sup>24</sup> if this is the case. While we typically use an alphabet of prime length, we may not always do so.

```
> AffineCode:= proc(plaintext::string, a::integer, b::integer)
  local textnum, codenum, i, p;
  global Alphabet;

  p      := length(Alphabet);
  if (gcd(a,p)>1) then
    error("The %-1 parameter %2 must be prime to the length of Alphabet %3",
          2,a,p);
  fi;
  textnum := StringToList(plaintext);
  codenum := [seq( modp(a*textnum[i]+b, p), i=1..length(plaintext)) ];
  ListToString(codenum);
end:

> mess:=AffineCode("A fine mess, Stanley!",10,5);
      mess := “^7Lj;B71BmmN7Pw{;’BHA”
```

What about decoding a message encoded by an affine cipher? If we know the original `a` and `b`, we could either write `AffineDecode`, which calculates `modp( (textnum[i]-b)/a, p)` instead, or, we could use `AffineCode(1/a, -b/a)`, which is really the same thing. We’ll take the latter approach.

If you think about it for a second, you may wonder whether  $1/a \bmod p$  makes any sense at all. It does, but only if we interpret division in the right way. We should not divide by  $a$  in the “ordinary way” and then reduce modulo  $p$ , but instead interpret it to mean the solution  $x$  to the equation  $ax \equiv 1 \pmod{p}$ . That is,  $1/a \bmod p$  is the multiplicative inverse of  $a$  in  $\{0, 1, \dots, p-1\}$ . Such an element exists exactly when  $a$  and  $p$  are relatively prime, as a consequence of Prop. 7.1. Since Maple knows about division modulo  $p$ , it allows us to use this notation, which is quite convenient.

```
> 1/6 mod 97;
```

In our definition of `AffineCode` we insisted that `a` and `b` be integers. This means that using a call like `AffineCode(mess,1/10,-5/10)` will give an error, since  $1/10$  is not an integer.

---

<sup>24</sup>Versions prior to Maple 6 need to use `ERROR` instead, which has slightly different behavior.

We can fix that by changing the definition of `AffineCode` to be `proc(plaintext::string, a::rational, b::rational)` or we can explicitly calculate these modular expressions before passing them to `AffineCode`.

```
> AffineDecode:= proc(ciphertext::string, a::integer, b::integer)
  local A, B, p;
  global Alphabet;
  p := length(Alphabet);
  A := 1/a mod p;
  B := -b/a mod p;
  AffineCode(ciphertext, A, B);
end:
> AffineDecode(mess,10,5);
```

“A fine mess, Stanley!”

### 7.3 Breaking an affine cipher

Suppose we have a message that was encrypted with an affine cipher, and we also know the `Alphabet` in use. Suppose also that we know (or manage to guess) the encoding of two letters. How do we then decipher the message? It is quite straightforward: since the message was enciphered with  $x \mapsto y = (ax + b) \bmod p$ , and we have two examples of  $x$  and  $y$ , we need only solve a pair of linear equations mod  $p$ . This is quite straightforward to do by hand, but Maple will also do it for us.

For example, suppose we receive the message

nrIbQyxEbkrI bwr EUbJaPybL abpyJJy UbUlrrxU

which we know is encoded in the 53 letter alphabet consisting of the letters A–Z, a blank, then a–z. Suppose we also guess that the letter “b” (28) is the encoding for a space (26), and that “U” (20) is an encoded “s” (45). We can figure out the encoding transformation by solving the pair of equations

$$26a + b \equiv 28 \pmod{53} \quad 45a + b \equiv 20 \pmod{53}.$$

We have Maple do this for us using `msolve`.

```
> msolve({a*26 + b = 28, a*45 + b = 20},53);
```

$$\{b = 25, a = 47\}$$

Thus, we see that the message was encoded by  $x \mapsto 47x + 25 \pmod{53}$ , and can be decrypted using  $a = 1/47 \pmod{53}$  and  $b = -25/47 \pmod{53}$ . Of course, it would have been more efficient to solve for the *deciphering* transformation directly.

```
> msolve({A*28 + B = 26, A*20 + B = 45},53);
```

$$\{A = 44, B = 13\}$$

```
> Alphabet:="ABCDEFGHIJKLMNOPQRSTUVWXYZ abcdefghijklmnopqrstuvwxyz":
AffineCode("nrIbQyxEbkrI bwr EUbJaPybL abpyJJy UbUlrrxU", 44, 13);
```

“You bend your words like Uri Gellers spoons”

## 8 Enciphering matrices

### 8.1 Treating text as vectors

We can let our enciphering transformation  $f$  operate on larger blocks of text than just a single character at a time. However, instead of treating a multi-byte string as a “super-character”, as in §5.3, we could think of it as a vector of characters. Thus, if we were using, say, 2-vectors, then the name “Dougal” would become the three vectors [68, 111], [117, 103] and [97, 108] (assuming we are using the ASCII codes). If the length of the text we were converting was not divisible by 2, we would have to extend it by adding a character to the end. Typically, this would be a space, or perhaps a null character (ASCII 0). Thus the name “Ian” occurring at the end of a message would become “Ian ”, and might encode as the pair of vectors [73, 97] and [110, 32].

Below is an implementation of this conversion, `StringToVects`, which takes as input the text to convert to  $n$ -vectors and the length of the vectors,  $n$ . It uses the `StringToList` procedure from §4, as well as the usual Maple commands.

```
> StringToVects:= proc(text::string,n::posint)
    local i,j,textnums,vectors;
    textnums := StringToList(cat(text, seq(" ", i=1..modp(n-length(text),n))));
    [seq([seq(textnums[(j-1)*n + i], i=1..n)],
        j=1..nops(textnums)/n)];
end:
```

The first thing it does is extend the length of the text to a multiple of  $n$ , if needed. This relies on the observation that  $k + ((n - k) \bmod n)$  is a multiple of  $n$ , and on the fact that Maple will do the right thing with a command like `seq(i,i=1..0)`, that is, produce nothing whatsoever. Then the text, extended by the correct number of blanks, is converted to a list of numbers.

We now parcel up the list of numbers `textnums` into chunks of size  $n$ . This is done by noting that the characters in the  $j$ -th vector will come from positions  $n(j - 1) + 1$  through  $n(j - 1) + n$ ,

so the inner `seq` gives us the  $j$ -th vector, as the outer one ranges over all the vectors. If you don't see how this works right off, just take a moment to write out an explicit example for yourself.

Strictly speaking, the result is a list of lists, rather than a list of vectors. Maple is happy to accept a list in contexts where a vector is wanted.

Converting back from a list of  $n$ -vectors to a character string is easy, because we already have `ListToString` to convert a list of numbers to a character string for us. All that is needed is remove the inner structure from each of our vectors, leaving only the sequence of numbers. We can use `map` and `op` to do this for us.

```
> map(op,[[36, 79], [85, 71], [65, 76]]);
```

```
[36, 79, 85, 71, 65, 76]
```

One small problem is that for Maple, there is a difference between a vector and a list of numbers—a vector has a little more structure. While we can use the two almost interchangeably (in the same way that we can think of a point in  $\mathbb{R}^n$  as an  $n$ -vector, or of as an  $n$ -vector as a point), some functions such as `op` that deal with the internal structure will give different results.

```
> l := [vector([36,79]), vector([85, 71]), vector([65, 76])]:
map(op,l);
```

```
[1..2, [1 = 36, 2 = 79], 1..2, [1 = 85, 2 = 71], 1..2, [1 = 65, 2 = 76]]
```

In order to be sure that our function will work with both vectors and lists, we need to convert the vectors to lists first. This is easily done with the Swiss army knife `convert`. The unusual-looking type declaration in the function says that we will be happy to accept a list of lists, a list of vectors, or a list of `Vectors`,<sup>25</sup> provided their elements are integers.

```
> VectsToString:=
  proc(vects::list({list(integer), vector(integer), Vector(integer)}))
    local l;
    l:=map(convert, vects, list);
    ListToString(map(op, l)):
  end;
```

## 8.2 Affine encoding with matrices

Now that we know how to think of character strings as vectors, we can manipulate them. Our enciphering transformation is just a permutation of the  $n$ -vectors.

For example, the Vignère cipher (§5.1) can be interpreted as the applying the map  $v \mapsto v + b \pmod{p}$ , where  $v$  and  $b$  are vectors the length of the key.

---

<sup>25</sup>Maple has two different linear algebra packages, `linalg` and the newer `LinearAlgebra` introduced in Maple 6. The `vector` type is from the former, `Vector` from the latter.

We can do a more general trick using the affine transformation

$$v \mapsto Av + b \pmod{p}$$

where  $A$  is an  $n \times n$  matrix, and  $b$  and  $v$  are  $n$ -vectors. As before, we need certain restrictions on the matrix  $A$ . In particular, we need it to be an invertible matrix (so that we can decipher the encoding). This is guaranteed by requiring that the determinant of  $A$  be nonzero and share no common factors with  $p$ . As before, choosing an alphabet of prime length makes our lives much easier. The linear version of this encoding,  $v \mapsto Av$ , is sometimes called a Hill cipher after Lester Hill, who studied such cryptosystems in the late 1920s and early 1930s [Hil1], [Hil2].

In Maple, one does matrix arithmetic pretty much as you would expect, except that we need to use `&*` to indicate matrix multiplication. This is because matrix multiplication is not commutative, while `*` is. Using `&*` tells Maple not to attempt the multiplication right away, and so it doesn't carry out the operations until you execute `evalm`. When the `evalm` function is executed, it does the matrix multiplication properly. Thus, to compute  $Av + b$ , we would write `evalm(A &* v + b)`. Because we want to reduce everything modulo  $p$ , we map `modp` onto the matrix. Since we might want to specify an element as the inverse of another, we allow our matrix and vector to contain rationals. As before, we accept several types of inputs.

```
> with(linalg):
AffineMatEncode :=
proc(text::string,
    A::{matrix(rational),Matrix(rational),list(list(rational))},
    B::{vector(rational),Vector(rational),list(rational)})
local vtext,vcrypt,i,n,p;
global Alphabet;

p:=length(Alphabet);
n:=nops(convert(B,list)); # look at B to get the size

if ( gcd(linalg[det](A),p)>1) then
    error("Determinant of %1 is not prime to the length of Alphabet %2",
        A,p);
fi;

vtext := StringToVects(text, n);
vcrypt:= [seq( map(modp,
                    evalm( A &* vtext[i] + B),
                    p),
                i=1..nops(vtext))];
VectsToString(vcrypt);
end;
```

Note that matrices allow the encodings of nearby characters in the message to interact. For example, we can exchange every other character in a message using the matrix  $\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ :

We should point out that all the examples in this section are using the 53-character alphabet of the previous section.

```
> AffineMatEncode("I am so mixed up", [[0,1],[1,0]], [0,0]);
```

“ Imas oimex dpu”

To decode a message encoded using  $A, b$ , we can encode the crypttext with the matrix  $A^{-1}$  and the shift vector  $-A^{-1}b$ . You should check for yourself why this is the proper transformation. To obtain the inverse matrix in Maple, you can use `inverse`. This is part of the linear algebra package, so you must either first issue the command `with(linalg,inverse)` or refer to the function as `linalg[inverse]`.

```
> A := matrix([[12, 27], [46, 44]]); B:=[25, 50];
```

$$A := \begin{bmatrix} 12 & 27 \\ 46 & 44 \end{bmatrix}$$

$$B := [25, 50]$$

```
> AffineMatEncode("Lester Hill",A, B);
```

xMPKFUGRxtaa

```
> AffineMatEncode(%, inverse(A), evalm(-inverse(A) &* B));
```

“Lester Hill ”

### 8.3 A Known-plaintext attack on an affine matrix cipher

As in §7.3, we can determine the key if we know what certain characters correspond to. If the message is encoded using  $n$ -vectors, we will need to know the decodings for  $n + 1$  different vectors, that is,  $n(n + 1)$  characters. We then solve the  $n + 1$  vector equations simultaneously to find out the deciphering matrix and shift vector.

We will give a small example: Suppose we receive a message that we know was encoded using an affine cipher on 2-vectors, in our usual 97-letter alphabet. Suppose we also know that the plaintext “secret” corresponds to the ciphertext “cilbup”. This is sufficient information to decipher the message, because we know the encodings of three different 2-vectors.

```
> Alphabet := cat("\n\t", Select(IsPrintable,convert([seq(i,i=1..255)],bytes))):
plain:=StringToVects("secret",2);
```

$$plain := [[85, 71], [69, 84], [71, 86]]$$



```
> crypt:=StringToVects("cilbup",2);
      crypt := [[69, 75], [78, 68], [87, 82]]
```

These three vectors give us three equations to solve:

$$Ac_1 + B \equiv p_1 \quad Ac_2 + B \equiv p_2 \quad Ac_3 + B \equiv p_3$$

where the  $c_i$  are the crypttext vectors and the  $p_i$  are the corresponding plaintext, and the equations are taken mod 97. (If this seems backwards, remember we are trying to find the deciphering transformation.) If the code were linear, then the deciphering matrix would be  $PC^{-1}$ , where  $P$  and  $C$  are the matrices whose columns are made up of the vectors  $p_i$  and  $c_i$  respectively.<sup>26</sup>

Since the code is affine, we can either play around manipulating matrices, or we can just translate everything into 6 linear equations in 6 variables. We will take the latter approach.

```
> A:=matrix([[a,b],[c,d]]);
      B:=[e,f];
```

$$A := \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

$$B := [e, f]$$

Note that we wish to solve  $Ac_i + B \equiv p_i \pmod{97}$  for  $A$  and  $B$ . If we execute the statement `evalm(A &* crypt[i] + B - plain[i])`, we will get a vector which should be the zero vector. For example, using the first sample, we obtain

```
> evalm( A&*crypt[1] + B - plain[1]);
```

$$[69a + 75b - 85 + e, 69c + 75d - 71 + f]$$

This gives us a vector of two expressions we would like to set to zero. If we do this for all three of our samples, and unwrap the vectors using `op` and `convert( ,list)`, we will get a system of equations appropriate to solve:

```
> eqns:= {seq(op(convert(
      evalm(A&*crypt[i] + B - plain[i]),
      list)), i=1..3)};

      eqns := {87c + 82d - 86 + f, 87a + 82b - 71 + e, 78c + 68d - 84 + f,
      78a + 68b - 69 + e, 69c + 75d - 71 + f, 69a + 75b - 85 + e}
```

Now, we use `msolve` to find out the decoding matrix and shift vector.

---

<sup>26</sup>Do you see why this is so? Hint: think about the relationship between matrices and linear functions.

```
> Answer := subs(msolve(eqns,97),[evalm(A),B]);
```

$$Answer := \left[ \begin{bmatrix} 42 & 84 \\ 19 & 78 \end{bmatrix}, [5, 88] \right]$$

We can decode our message now, by encoding the crypttext using the matrix  $A$  and shift  $B$  found above.  $A$  is the inverse of the matrix used for encoding, and  $B$  is  $A$  times the negative of the original shift vector.

```
> AffineMatEncode("cilbup",Answer[1],Answer[2]);
```

“secret”

It is worth noting that combining an affine matrix cipher with a large multibyte alphabet (see §5.3) greatly increases its security, at least if one restricts to short messages, each using different keys. While the matrix cipher is still susceptible to a known-plaintext attack, getting the plaintext becomes much harder.

## 9 Modern cryptography

### 9.1 Secure cryptosystems

As mentioned before, most of the ciphers we have examined are not really cryptographically secure, although they may have been adequate prior to the widespread availability of computers. A cryptosystem is called **secure** if a good cryptanalyst, armed with knowledge the details of the cipher (but not the key used), would require a prohibitively large amount of computation to decipher a plaintext. This idea (that the potential cryptanalyst knows everything but the key) is called **Kerckhoff’s Law** or sometimes **Shannon’s Maxim**.

Kerckhoff’s Law at first seems unreasonably strong rule; you may be thinking that, given a mass of encrypted data, how is the cryptanalyst to know *by what means* it was encrypted? If you are encrypting your own files for personal use and writing your own software which you keep under lock and key, this assumption is not unreasonable. But today, most encryption is done by software or hardware that the user did not produce himself. One can reverse engineer a piece of software (or hardware) and make the cryptographic algorithm apparent, so we cannot rely on the secrecy of the method alone as a good measure of security. For example, when you make a purchase over the internet, the encryption method that is used between your browser and the seller’s web server is public knowledge. Indeed, it must be public: if not, only those privy to the secret could create web browsers or web servers. The security of the transaction depends on the security of the keys.

There are several widely used ciphers which are believed to be fairly secure. We will not go into the details of any of these; most of the algorithms are quite involved and require manipulations at the bit level— Maple is not really the appropriate tool for them. Implementations

are widely available, both commercially and as free public software (see, for example, the International Cryptography page<sup>27</sup>).

Probably the most commonly used cipher today is DES (the Data Encryption Standard), which was developed in the 1970s and has been adopted as a standard by the US government. The standard implementation of DES operates on 64-bit blocks (that is, it uses an alphabet of length  $2^{64}$ — each “character” is 8 bytes long), and uses a 56-bit key. Unfortunately, the 56-bit key means that DES, while secure enough to thwart the casual cryptanalyst, is attackable with special hardware by governments, major corporations, and probably well-heeled criminal organizations. One merely needs to try a large number of the possible  $2^{56}$  keys.<sup>28</sup> This is a formidable, but not insurmountable, computational effort. A common variation of DES, called Triple-DES, uses three rounds of regular DES with three different keys (so the key length is effectively 168 bits), and is considerably more secure.

DES was originally expected to be used for “only a few years” when it was first designed. However, due to its surprising resistance to attack it was generally unassailable for nearly 25 years. The powerful methods of linear and differential cryptanalysis were developed to attack block ciphers like DES.

In the January of 1997, the National Institute for Standards and Technology (NIST) issued a call for a new encryption standard to be developed, called AES (the Advanced Encryption Standard). The requirements were that they operate on 128-bit blocks and support key sizes of 128, 192, and 256 bits. There were five ciphers which advanced to the second round: MARS, RC6, Rijndael, Serpent, and Twofish. All five were stated to have “adequate security”—Rijndael was adopted as the standard in October of 2000. Rijndael was created by the Belgian cryptographers Joan Daemen and Vincent Rijmen; the underlying mathematics is the algebra of polynomials in the finite field  $GF(2)$ .

Other commonly used ciphers are IDEA (the International Data Encryption Algorithm, developed at the ETH Zurich in Switzerland), which uses 128-bit keys, and Blowfish (developed by Bruce Schneier), which uses variable length keys of up to 448 bits. Both of these are currently believed to be secure. Another common cipher, RC4 (developed by RSA Data Security) can use variable length keys and, with sufficiently long keys, is believed to be secure. Some versions the Netscape web browser used RC4 with 40 bit keys for secure communications. A single encrypted session was broken in early 1995 in about 8 days using 112 networked computers (see Damien Doligez’s web page about cracking SSL<sup>29</sup> for more details); later in the same year a second session was broken in under 32 hours. Given the speed increases in computing since then, it is reasonable to believe that a 40-bit key can be cracked in a few hours. Notice, however, that both of these attacks were essentially brute-force, trying a large fraction of the  $2^{40}$  possible keys. Increasing the key size resolves that problem. Nearly all browsers these days use at least 128-bit keys.

---

<sup>27</sup><http://www.cs.hut.fi/ssh/crypto/>

<sup>28</sup>that is, 72,057,594,037,927,936 keys.

<sup>29</sup><http://pauillac.inria.fr/doligez/ssl/>

There are, of course, many other ciphers in common use, and more being developed all the time.

## 9.2 Message digests

Related to encryption are message digests, or cryptographic hash functions. A message digest takes an arbitrary length string and computes a number, or hash, from it. Changing the string changes the value of the hash. Hashing functions also play a major role in searching and sorting, and hence in computer databases. While it is obvious that several plaintexts must correspond to the same hash value, it is a difficult problem in general to construct such a plaintext from the value of the hash alone. Thus, one can use a good message digest to **authenticate** a message—appending an encrypted version of the hash to an otherwise clear text message can act as a seal, showing that the message was not otherwise tampered with, since it would be infeasibly hard to construct a different message with the same hash.

Some commonly used Message digests are MD5 (developed by RSA Data Security), and SHA or SHS (the Secure Hash Algorithm/Standard, developed by the US government).

It is worth remarking here that computer systems can use a cryptographic hash function as a way to store a password (for example, Linux allows the use of MD5 for passwords). A computer password never need be decrypted; it only needs to be verified. Instead of storing the password itself, or an encrypted version of it, the computer stores the corresponding hash. Then, when a user signs in, the password is hashed and compared to the stored value. If they match, the user typed the right password (or managed to find another password with the same hash, which is computationally impossible).

## 9.3 Public Key cryptography

All of the ciphers so far discussed are symmetric encryption routines (also called secret-key or traditional ciphers): if one knows the key used to encode the message, one also can decrypt the message without too much work. For example, in an affine cipher  $x \mapsto ax + b \pmod p$  (see §7), if we know the values of  $a$  and  $b$  used to encode the message, we can quickly compute the deciphering key  $1/a \pmod p$ ,  $-b/a \pmod p$ . By the same token, the ability to decrypt the message usually also yields the key used for encryption.<sup>30</sup>

---

<sup>30</sup>A story related to this appears in Casanova's memoirs (1757) (quoted in [Kahn]). He was given an encrypted manuscript by a Madame d'Urfé, which he deciphered. From deciphering, he determined the key used to encipher the manuscript. Later, upon learning this, Madame d'Urfé was incredulous, believing he must have learned the keyword in order to decode the message.

I then told her the key-word, which belonged to no language, and I saw her surprise. She told me that it was impossible, for she believed herself the only possessor of that word which she kept in her memory and which she had never written down.

I could have told her the truth—that the same calculation which had served me for deciphering the manuscript had enabled me to learn the word—but on a caprice it struck me to tell her that

This was true of all cryptosystems up until the mid 1970s: knowledge of how to encode a message allowed one also to decipher it. However, in 1976 Whitfield Diffie and Martin Hellman invented **public key cryptography** [DH76]. In a public key system, someone who knows how to encipher a message cannot determine how to decipher the message without a prohibitively large computation. That is, given the enciphering transformation  $f_{key} : \mathcal{P} \mapsto \mathcal{C}$ , discovering  $f_{key}^{-1}$  is not computationally realistic. It is important to realize that there may be a procedure for doing so, but to carry out this process would take a prohibitively long time (say, hundreds of years on the fastest known computers). Such a function is sometimes called a **one-way function**. However, to be usable for public key cryptography, it must be possible to find the inverse of the function provided you have some extra information. This extra information is referred to as a **trapdoor** (like a trapdoor that a magician would use to escape from a “sealed” box), resulting in a **trapdoor one-way function**.

What good is a public key cryptosystem? Since knowledge of the encoding key  $K_E$  does not give information about the decoding key  $K_D$ , the encoding key  $K_E$  can be made public (hence the name public key). This allows anyone to encode messages that only the recipient can decode. For example, suppose you want to make a credit card purchase over the Internet. Data sent across the Internet unencrypted is not secure: someone with access to the transmission lines can listen in and capture sensitive data.<sup>31</sup> However, if the merchant provides a public key, you can encrypt your credit card number and transmit it without worry. Only the merchant can decrypt the message, even though anyone may send one.<sup>32</sup>

Another feature of public key cryptography is **authentication**. In many public key systems, either the enciphering key  $K_E$  or the deciphering key  $K_D$  may be made public without revealing the other one. In order to digitally sign a message, I could append my name encrypted with my enciphering key, which I keep secret. I could also make public the deciphering key  $K_D$ . Then anyone can check that I was the actual sender, because only I could have encoded my name. Of course, to stop someone from merely appending a valid encrypted signature to a forged message, I should also encrypt something specific to that message, such as a message digest (see §9.2) computed on its contents. This feature is commonly used to digitally sign email by

---

a genie had revealed it to me. This false disclosure fettered Madame d’Urfé to me. That day, I became the master of her soul, and I abused my power. Every time I think of it, I am distressed and ashamed, and I do penance now in the obligation under which I place myself of telling the truth in writing my memoirs.

<sup>31</sup>This is not as hard as you might expect. In many places, such as universities, portions of the network backbone pass through public places like dormitories or classrooms. One merely needs to attach a computer to the network lines to “sniff” the network traffic. This will, of course, be limited to users communicating to or from that location, but that can be quite a lot. This same problem occurs with some cable-modem providers: appropriate software can detect the network traffic of one’s neighbors. Also, there are many segments of the Internet which consist of microwave links. To listen in on a huge part of network traffic, all you need is an appropriately placed receiver.

<sup>32</sup>Of course, you should still only deal with reputable merchants. You probably wouldn’t give your credit card number to a guy selling stuff on the street; neither should you automatically trust anyone with a web site.

software such as PGP<sup>33</sup>.

Public-key cryptography is very computationally intensive, so typically its use is limited to allow for the secure transmission of a secret key; this secret key is then used to encrypt the rest of the message using a symmetric encryption method such as AES, triple-DES or IDEA.

A common public key system in current use is RSA (named for its inventors: Ronald Rivest, Adi Shamir and Leonard Adelman), which we will look at in detail in §11. Others are the Diffie-Hellman key exchange system, El Gamal, Subset-sum (or Knapsack) systems, and systems based on elliptic curves. We will not discuss these here.

## 10 Some Number Theory

Most public key systems rely on number-theoretic results. Before we can discuss the implementation of one, we need to quickly go over the necessary background. We have already used a tiny amount of number theory (in our discussion of computing  $\text{mod } p$  and of the greatest common divisor). Of course, this must be done briefly, and we will only touch on a small part of a large and ancient field—the interested reader would do well to consult a text on number theory (e.g. [NZM], [Ros]) for more information.

### 10.1 The greatest common divisor and the Euclidean algorithm

We have already met the greatest common divisor, or gcd, which is the largest integer which divides both of a pair of numbers. Two numbers are said to be **relatively prime** if their greatest common divisor is 1. As we have already seen, finding two relatively prime numbers has important applications in many cryptosystems.

How can we determine the gcd of two numbers? If the numbers are not too large, just looking at their factors does the trick. For example,

$$\text{gcd}(138, 126) = 6 \quad \text{since} \quad 138 = 2 \cdot 3 \cdot 23 \quad \text{and} \quad 126 = 2 \cdot 3^2 \cdot 7.$$

However, there must be a more efficient way, since Maple can calculate the gcd of two 60-digit numbers in well under a second, while taking a long time (about 10 minutes on a 600 Mhz pentium computer) to factor either one of them.

One such algorithm is the Euclidean algorithm, which has been around for thousands of years. It works by computing successive differences—to compute  $\text{gcd}(a, b)$ , where  $a > b$ , we first compute  $r_1 = a - k_1b$ , where  $k_1b$  is the largest multiple of  $b$  that is less than  $a$ . Then we repeat this, finding  $r_2 = b - k_2r_1$ ,  $r_3 = r_1 - k_3r_2$ , and so on until the remainder is zero. The last nonzero remainder is the gcd. For example, to calculate  $\text{gcd}(138, 126)$ ,

---

<sup>33</sup><http://www.pgpi.com>

$$\begin{array}{rclcl}
 a & & = & 138 & \\
 & & & b & = 126 \\
 & & & a & - b = 12 = r_1 \\
 b - 10r_1 & = & -10a + 11b & = & 6 = r_2 \\
 r_1 - 2r_2 & = & 21a - 23b & = & 0
 \end{array}$$

This not only gives  $\gcd(138, 126) = 6$ , but expresses it as a difference of multiples of the two numbers:  $6 = 11 \times 126 - 10 \times 138$ . We can write this procedure more formally as a theorem.

**Theorem 10.1** (The Euclidean Algorithm) *Let  $a$  and  $b$  be two positive integers, with  $b < a$ . Then we can construct two decreasing sequences of positive integers  $r_i$  and  $k_i$  as follows:*

$$\begin{array}{rclcl}
 a & - & k_1 b & = & r_1 & & 0 < r_1 < b \\
 b & - & k_2 r_1 & = & r_2 & & 0 < r_2 < r_1 \\
 r_1 & - & k_3 r_2 & = & r_3 & & 0 < r_3 < r_2 \\
 r_2 & - & k_4 r_3 & = & r_4 & & 0 < r_4 < r_3 \\
 & & & & \vdots & & \\
 r_{n-2} & - & k_n r_{n-1} & = & r_n & & 0 < r_n < r_{n-1} \\
 r_{n-1} & - & k_{n+1} r_n & = & 0 & & 
 \end{array}$$

and  $r_n$  is the greatest common divisor of  $a$  and  $b$ .

We have already seen that the Euclidean algorithm also gives us the following:

**Corollary 10.2** *Let  $a$  and  $b$  be two positive integers. Then there are integers  $x$  and  $y$  so that  $ax + by = \gcd(a, b)$ .*

Why does the Euclidean Algorithm work? It follows from the observation that if  $a = bk + r$ , then  $\gcd(a, b) = \gcd(b, r)$ . For example,  $\gcd(138, 126) = \gcd(126, 12) = \gcd(12, 6) = 6$ . We write this as a lemma.

**Lemma 10.3** *Let  $b$  be a positive integer, and  $a$  be a nonnegative integer. If*

$$a = bk + r, \quad \text{with } k \text{ and } r \text{ positive integers}$$

then  $\gcd(a, b) = \gcd(b, r)$ .

**Proof.** Let  $d = \gcd(a, b)$ . Then since  $d$  is a divisor of  $a$  and a divisor of  $b$ , it is also a divisor of  $a - bk$ . Since  $r = a - bk$ , we have  $d$  a divisor of  $r$ . So it must divide  $\gcd(b, r)$ .

Now, in order to prove equality, we want to show that also  $\gcd(b, r)$  is a divisor of  $d$ . Certainly  $\gcd(b, r)$  divides  $a$ , since  $a = r + bk$ . It also divides  $b$ , and so we have  $\gcd(b, r)$  as a divisor of  $\gcd(a, b)$ .

Thus,  $\gcd(b, r) = \gcd(a, b)$ . □

Now we are ready to prove that the Euclidean Algorithm (Thm. 10.1) always yields the greatest common divisor.

**Proof.** Notice that the sequence  $b > r_1 > r_2 > r_3 > \dots$  is a strictly decreasing sequence of positive integers, so it must eventually terminate with some  $r_n > 0$ . The last equation can be written as  $k_{n+1}r_n = r_{n-1}$ , and so  $r_n$  is a divisor of  $r_{n-1}$ . Then certainly  $r_n = \gcd(r_n, r_{n-1})$ . Applying the lemma to the equation above it, we obtain  $\gcd(r_n, r_{n-1}) = \gcd(r_{n-1}, r_{n-2})$ . Repeating in this way, we obtain  $r_n = \gcd(r_n, r_{n-1}) = \gcd(r_{n-1}, r_{n-2}) = \dots = \gcd(r_1, b) = \gcd(a, b)$ . □

Recall that if  $r$  is such that  $ar \equiv 1 \pmod{n}$ , then  $r$  is called the multiplicative inverse of  $a$  modulo  $n$ . We have already dealt with such objects extensively. We now state a theorem which tells us exactly when such inverses exist. Note that this is really just Prop. 7.1 slightly reworded.

**Theorem 10.4** *Let  $a$  and  $n$  be integers, with  $n \geq 2$ . Then  $a$  has a multiplicative inverse modulo  $n$  if and only if  $\gcd(a, n) = 1$ .*

**Proof.** First, suppose that  $a$  has an inverse modulo  $n$ . Then there is a  $k$  so that

$$ak \equiv 1 \pmod{n}.$$

In particular,  $n$  is a divisor of  $ak - 1$ , so there must be some integer  $t$  so that  $ak - 1 = nt$ . Rewriting this as

$$ak - nt = 1,$$

we see that  $\gcd(a, n) = 1$  (by Cor. 10.2).

Now suppose that  $\gcd(a, n) = 1$ . Then again using Cor. 10.2, there are integers  $r$  and  $s$  with  $ar + ns = 1$ . This means  $ar - 1$  is divisible by  $n$ , so

$$ar \equiv 1 \pmod{n}.$$

But this says that  $r$  is the inverse of  $a$  modulo  $n$ . □

**Notation.** It is common (and quite convenient) to denote the set of integers modulo  $n$  as  $\mathbb{Z}_n$ . That is,

$$\mathbb{Z}_n = \{0, 1, 2, \dots, n - 1\}$$

As we already know, this set is closed under addition and multiplication, and the usual associative, distributive, and commutativity laws hold. This set is an example of an algebraic structure called a ring.



**Notation.** We will use  $\mathbb{Z}_n^*$  to denote the elements of  $\mathbb{Z}_n$  which have multiplicative inverses. In light of the last theorem, we have

$$\mathbb{Z}_n^* = \{a \in \mathbb{Z}_n \mid \gcd(a, n) = 1\}$$

For example,

$$\mathbb{Z}_7^* = \{1, 2, 3, 4, 5, 6\} \quad \mathbb{Z}_8^* = \{1, 3, 5, 7\} \quad \mathbb{Z}_{12}^* = \{1, 5, 7, 11\}$$

## 10.2 The Chinese Remainder Theorem

In this section, we state a very old theorem about simultaneous solutions to congruences. This theorem may have been known to the eight-century Buddhist monk I-Hsing, and certainly appears in Ch'in Chiu-shao's *Mathematical Treatise in Nine Sections*, written in 1247. It is sometimes (rarely) referred to as the "Formosa Theorem".

**Theorem 10.5** *Let  $m \geq 2$  and  $n \geq 2$  be integers with  $\gcd(m, n) = 1$ . Then for any integers  $a$  and  $b$ , the system of equations*

$$\begin{aligned} x &\equiv a \pmod{m} \\ x &\equiv b \pmod{n} \end{aligned}$$

*has a solution. Furthermore, there is only one solution  $x$  between 0 and  $mn - 1$ .*

As an example of this, suppose five kids pool all their money and buy a bag of candy. They divide the candy up evenly, and find there are 3 pieces left over. Suddenly, one kid's mom comes out, tells him "no candy before dinner", and hauls him off. The other four kids divide up the candy and there is one piece left over. While they are squabbling over it, a dog runs up and eats the disputed piece.

The Chinese Remainder Theorem tells us that we can always determine how many pieces of candy were in the bag, at least up to multiples of 20. In the above example, there were 13 pieces of candy (or 33, or 53, or ...). If we want to solve this problem with Maple, the command `chrem([3,1],[5,4])` does the trick. We could also use `numtheory[mcombine](5,3,4,1)`.

**Proof.** Since  $m$  and  $n$  are relatively prime, there are integers  $k$  and  $t$  so that  $mk + nt = 1$ . Then

$$c = bkm + ant$$

is the desired solution. To verify this, we need only check:

$$c \equiv ant \equiv a(1 - mk) \equiv a \pmod{m}$$

and

$$c \equiv bkm \equiv b(1 - nt) \equiv b \pmod{n}$$

To check uniqueness, suppose there are two solutions  $c$  and  $d$ . Since  $c \equiv a \pmod{m}$  and  $d \equiv a \pmod{m}$ , we have  $c - d \equiv 0 \pmod{m}$ . Similarly,  $c - d \equiv 0 \pmod{n}$ . Since both  $m$  and  $n$  divide  $c - d$ , and  $m$  and  $n$  are relatively prime, we have that  $mn$  divides  $c - d$ . Thus  $c - d \equiv 0 \pmod{mn}$ , that is,  $c \equiv d \pmod{mn}$ .  $\square$

### 10.3 Powers modulo $n$

**Notation.** In this and future sections, we shall use the notation  $\llbracket a \rrbracket_b$  to denote  $a \pmod{b}$ .

First, let's take a look at a few examples of what happens when we reduce the sequence  $a, a^2, a^3, a^4, \dots$  modulo  $n$ . Consider  $\llbracket 3^k \rrbracket_{20}$ :

$$3, \llbracket 3^2 \rrbracket_{20} = 9, \llbracket 3^3 \rrbracket_{20} = 7, \llbracket 3^4 \rrbracket_{20} = \llbracket 21 \rrbracket_{20} = 1, \llbracket 3^5 \rrbracket_{20} = 3, \dots$$

The sequence consists of four numbers 3, 9, 7, 1 and then it repeats. What about 4? We have

$$4, \llbracket 4^2 \rrbracket_{20} = 16, \llbracket 4^3 \rrbracket_{20} = 4, \llbracket 4^4 \rrbracket_{20} = 16, \dots$$

This is also periodic, but notice that we never get back to 1. Similarly, let's look at the powers of 6:

$$6, \llbracket 6^2 \rrbracket_{20} = \llbracket 36 \rrbracket_{20} = 16, \llbracket 6^3 \rrbracket_{20} = \llbracket 16 \cdot 6 \rrbracket_{20} = \llbracket 96 \rrbracket_{20} = 16, \llbracket 6^4 \rrbracket_{20} = 16, \dots$$

This gets stuck at 16. Finally, consider the powers of 11:

$$11, \llbracket 11^2 \rrbracket_{20} = \llbracket 121 \rrbracket_{20} = 1, \llbracket 11^3 \rrbracket_{20} = 11, \dots$$

This sequence is periodic of period 2, and contains 1.

**Definition 10.6** *If there is a positive integer  $k$  so that  $\llbracket a^k \rrbracket_n = 1$ , we say that  $a$  has finite multiplicative order modulo  $n$ . The smallest such integer  $k$  is called the order of  $a$  modulo  $n$ .*

We can readily categorize the elements with finite order. If you look at the examples above, 3 and 11 have finite order mod 20, while 4 and 6 do not. You probably have already guessed that this is because  $\gcd(3, 20) = 1$  and  $\gcd(11, 20) = 1$ , while  $\gcd(4, 20) = 4$  and  $\gcd(6, 20) = 2$ .

**Theorem 10.7** *Let  $a$  and  $n > 2$  be integers. Then  $a$  has finite order modulo  $n$  if and only if  $\gcd(a, n) = 1$ .*

**Proof.** First, suppose  $a$  has finite order mod  $n$ . Then there is some  $k$  so that

$$1 = \llbracket a^k \rrbracket_n = \llbracket a \rrbracket_n \cdot \llbracket a^{k-1} \rrbracket_n.$$

Thus,  $a$  has a multiplicative inverse, namely  $\llbracket a^{k-1} \rrbracket_n$ . By Thm. 10.4, this means  $\gcd(a, n) = 1$ .

Now we suppose  $\gcd(a, n) = 1$ , and we want to find the inverse of  $a$ . Consider the sequence

$$a, \llbracket a^2 \rrbracket_n, \llbracket a^3 \rrbracket_n, \dots, \llbracket a^n \rrbracket_n, \llbracket a^{n+1} \rrbracket_n.$$

Since there are  $n + 1$  elements and we are reducing modulo  $n$ , there must be at least two which are equal. Let's denote these by  $\llbracket a^k \rrbracket_n$  and  $\llbracket a^t \rrbracket_n$ , with  $1 \leq t < k$ . Since

$$\llbracket a^k \rrbracket_n = \llbracket a^t \rrbracket_n$$

we have

$$\llbracket a^{k-t} \rrbracket_n = 1.$$

Thus,  $a$  has order  $k - t \pmod n$ . □

So, we have seen that any invertible element of  $\mathbb{Z}_n$  has finite order. In particular, if  $p$  is a prime number, every element except 0 is invertible and so has finite order. What are the possible orders? A theorem of Pierre de Fermat gives a result in that direction.

**Theorem 10.8** (Fermat's Little Theorem, 1640) *Let  $p$  be a prime. If  $\gcd(a, p) = 1$ , then*

$$a^{p-1} \equiv 1 \pmod p.$$

*For any integer  $a$ ,*

$$a^p \equiv a \pmod p.$$

Among other things, this theorem says that if  $p$  is prime, then the order of  $a \pmod p$  always divides  $p - 1$ . Note that it doesn't say the order is  $p - 1$ , just that it is a divisor of it. Also, it says nothing about order with modulo a non-prime. We won't give the proof of this theorem here. Rather, we will give a proof of a more general result in the next section.

## 10.4 The Euler $\varphi$ -function and Euler's Theorem

Recall that we denote the set of invertible elements of  $\mathbb{Z}_n$  by  $\mathbb{Z}_n^*$ . That is,

$$\mathbb{Z}_n^* = \{a \in \mathbb{Z}_n \mid \gcd(a, n) = 1\}.$$

This set is sometimes called a **reduced residue system modulo  $n$** .

**Definition 10.9** *Let  $n$  be a positive integer. Then  $\varphi(n)$  is the number of elements in  $\mathbb{Z}_n^*$ , that is, the number of positive integers less than  $n$  which are relatively prime to  $n$ . This function is called the Euler  $\varphi$ -function,<sup>34</sup> or sometimes Euler's totient function.*

---

<sup>34</sup>Although the function was invented by Euler, the notation  $\varphi(n)$  to represent it is actually due to Gauss.

We can compute a few values of  $\varphi(n)$  just by counting.

$$\begin{array}{llll}
 \varphi(2) & = & 1 & \text{since } \mathbb{Z}_2^* = \{1\} \\
 \varphi(3) & = & 2 & \text{since } \mathbb{Z}_3^* = \{1, 2\} \\
 \varphi(4) & = & 2 & \text{since } \mathbb{Z}_4^* = \{1, 3\} \\
 \varphi(5) & = & 4 & \text{since } \mathbb{Z}_5^* = \{1, 2, 3, 4\} \\
 \varphi(6) & = & 2 & \text{since } \mathbb{Z}_6^* = \{1, 5\} \\
 \varphi(7) & = & 6 & \text{since } \mathbb{Z}_7^* = \{1, 2, 3, 4, 5, 6\} \\
 \varphi(8) & = & 4 & \text{since } \mathbb{Z}_8^* = \{1, 3, 5, 7\} \\
 \varphi(9) & = & 6 & \text{since } \mathbb{Z}_9^* = \{1, 2, 4, 5, 7, 8\} \\
 \varphi(10) & = & 4 & \text{since } \mathbb{Z}_{10}^* = \{1, 3, 7, 9\} \\
 \varphi(11) & = & 10 & \text{since } \mathbb{Z}_{11}^* = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\} \\
 \varphi(12) & = & 4 & \text{since } \mathbb{Z}_{12}^* = \{1, 5, 7, 11\}
 \end{array}$$

Of course, Maple also knows about the function  $\varphi$  in the `numtheory` package, so we could generate more values by asking Maple. But let's try to understand a bit more. We can readily prove some properties of  $\varphi$ .

**Proposition 10.10** *Let  $p \geq 2$  be a prime. Then*

- a.  $\varphi(p) = p - 1$ .
- b. If  $n$  is a positive integer, then  $\varphi(p^n) = p^n - p^{n-1}$ .
- c. If  $a$  and  $b$  are relatively prime, then  $\varphi(ab) = \varphi(a)\varphi(b)$ .

**Proof.** (a) This is immediate, since all non-zero elements of  $\{0, 1, \dots, p-1\}$  are invertible.

(b) We prove this part just by counting. Since  $p$  is prime, the only numbers in  $\mathbb{Z}_p$  which are not relatively prime to  $p^n$  are multiples of  $p$ . These are  $0, p, 2p, 3p, \dots, p^2, (p+1)p, \dots, (p^{n-1}-1)p$ . Note that  $p^{n-1}p = p^n$ , so there are exactly  $p^{n-1}$  multiples of  $p$  less than  $p^n$  (including 0). All the others are relatively prime to  $p$ —this gives the result.

(c) We need to show that the number of elements in  $\mathbb{Z}_{ab}^*$  is the same as the product of the size of  $\mathbb{Z}_a^*$  and  $\mathbb{Z}_b^*$ . We do that by showing that for each  $t \in \mathbb{Z}_{ab}^*$ , there is exactly one pair  $(r, s)$  with  $r \in \mathbb{Z}_a^*$  and  $s \in \mathbb{Z}_b^*$ .

So, let  $r \in \mathbb{Z}_a^*$  and  $s \in \mathbb{Z}_b^*$ . Then, by the Chinese Remainder Theorem (Thm. 10.5), there is an integer  $t < ab$  with

$$t \equiv r \pmod{a} \quad \text{and} \quad t \equiv s \pmod{b}.$$

Since  $t$  is a solution of the first, we can find an integer  $k$  so that  $r = t + ka$ , so  $\gcd(t, a) = \gcd(a, r) = 1$  by Lemma 10.3. Similarly,  $\gcd(t, b) = 1$ . Since  $t$  is relatively prime to both  $a$  and  $b$ , it is relatively prime to their product, so  $\gcd(t, ab) = 1$ . Thus  $t \in \mathbb{Z}_{ab}^*$ .

Now take  $t \in \mathbb{Z}_{ab}^*$ , and we must find a corresponding  $r$  and  $s$ . Let  $r = \llbracket t \rrbracket_a$ . Now, since  $\gcd(t, ab) = 1$ , certainly  $\gcd(t, a) = 1$ . Since  $r \equiv t \pmod{a}$ , there is an integer  $k$  so that  $t = r + ka$ , so we have  $\gcd(r, a) = 1$ . Thus  $r \in \mathbb{Z}_a^*$ . Similarly, we let  $s = \llbracket t \rrbracket_b$ , and we see that  $s \in \mathbb{Z}_b^*$ .

Since we have paired each  $t \in \mathbb{Z}_{ab}^*$  with a unique pair  $(r, s) \in \mathbb{Z}_a^* \times \mathbb{Z}_b^*$ , they have the same cardinality. The first set has  $\varphi(ab)$  elements, and the latter has  $\varphi(a)\varphi(b)$ .  $\square$

We now come to Euler's generalization of Fermat's result. This result is the central idea underlying the RSA public key cryptosystem.

**Theorem 10.11** (Euler, 1750) *Let  $a$  and  $n$  be relatively prime integers, with  $n \geq 2$ . Then*

$$\llbracket a^{\varphi(n)} \rrbracket_n = 1.$$

**Proof.** Let  $a$  be relatively prime to  $n$ , and consider the set

$$a\mathbb{Z}_n^* = \{ \llbracket ab \rrbracket_n \mid b \in \mathbb{Z}_n^* \}.$$

Since  $a$  and  $n$  are relatively prime,  $\llbracket a \rrbracket_n \in \mathbb{Z}_n^*$ . Thus every element of  $a\mathbb{Z}_n^*$  is in  $\mathbb{Z}_n^*$ . But also every element of  $\mathbb{Z}_n^*$  is an element of  $a\mathbb{Z}_n^*$ , since  $\llbracket ab \rrbracket_n = \llbracket ac \rrbracket_n$  if and only if  $\llbracket b \rrbracket_n = \llbracket c \rrbracket_n$ . This means that the sets  $a\mathbb{Z}_n^*$  and  $\mathbb{Z}_n^*$  are equal.

Now, multiply all the elements of  $\mathbb{Z}_n^*$  together; call the result  $N$ . Then  $\llbracket N \rrbracket_n \in \mathbb{Z}_n^*$ . If we multiply all the elements of  $a\mathbb{Z}_n^*$  together, we get  $a^{\varphi(n)}N$ , and since the two sets are the same, we have

$$\llbracket a^{\varphi(n)}N \rrbracket_n = \llbracket N \rrbracket_n.$$

Thus,

$$\llbracket a^{\varphi(n)} \rrbracket_n = 1.$$

$\square$

## 11 The RSA Public key cryptosystem

In 1978, R. Rivest, A. Shamir, and L. Adelman published the description of the RSA public key system [RSA], which we will describe here. As discussed in §9.3, a public key system allows anyone to encode messages that only the intended recipient may decode. This algorithm relies on the fact that unless one knows some special facts about  $n$ , calculating the Euler  $\varphi$ -function  $\varphi(n)$  is as hard as factoring  $n$ , which, for very large  $n$  (say, 200 digits), is very hard indeed.

## 11.1 Details of the RSA algorithm

To set up the system, we pick at random two large primes  $p$  and  $q$ , of about 100 digits each. We then compute the base  $n$  as the product of  $p$  and  $q$ . Note that  $\varphi(n) = (p-1)(q-1)$ . We also pick some other number  $e$ , called the **exponent**, which is relatively prime to  $\varphi(n)$ . We make the numbers  $(n, e)$  public—these form the key needed for encoding. We also use the Euclidean algorithm to compute  $x$  which satisfies  $ex - y\varphi(n) = 1$ ; that is, so that  $\llbracket ex \rrbracket_{\varphi(n)} = 1$ . The number  $x$  is the part of the key we keep private; our decrypting key is the pair  $(n, x)$ .

**To encode a message:** The sender divides his message up into blocks of length  $k$ , and to each block assigns an integer  $\beta$  (for example, by treating the characters of the message as a  $k$ -graph, as in §5.3, with  $0 \leq \beta < n$ ). For each block of plaintext, the sender transmits  $\llbracket \beta^e \rrbracket_n = \mu$ .

**To decode the message:** For each unit  $\mu$  of the crypttext received, the recipient computes  $\llbracket \mu^x \rrbracket_n$ , which is the integer representing of a  $k$ -graph of the original plaintext message. We can see that this is so, because

$$\llbracket \mu^x \rrbracket_n = \llbracket \beta^{ex} \rrbracket_n = \llbracket \beta^{1+\varphi(n)y} \rrbracket_n = \llbracket \beta \rrbracket_n \llbracket \beta^{y\varphi(n)} \rrbracket_n = \llbracket \beta \rrbracket_n (\llbracket \beta^{\varphi(n)} \rrbracket_n)^y = \llbracket \beta \rrbracket_n \cdot 1^y = \beta.$$

This relies on Euler's theorem (Thm. 10.11), so that  $\llbracket \beta^{\varphi(n)} \rrbracket_n = 1$ . Euler's theorem only guarantees that this will work for  $\gcd(\beta, n) = 1$ . However, we can easily see that decryption works for all  $\beta < n$ , whether or not it is relatively prime to  $n$ .

**Proposition 11.1** *Let  $n = pq$ , with  $p$  and  $q$  prime, and let  $0 \leq \beta < n$ . Then for any integer  $y$ ,  $\llbracket \beta^{1+y\varphi(n)} \rrbracket_n = \beta$ .*

**Proof.** If  $\gcd(\beta, n) = 1$ , the result follows from Euler's theorem.

Otherwise, since  $p$  and  $q$  are prime, we have  $\beta$  a multiple of either  $p$  or  $q$ , but not both (since  $\beta < pq$  by assumption). Suppose  $\beta \equiv 0 \pmod{p}$ , and so

$$\llbracket \beta^{1+y\varphi(n)} \rrbracket_p = 0 = \llbracket \beta \rrbracket_p.$$

Since  $\beta$  is a multiple of  $p$  less than  $pq$ , and both  $p$  and  $q$  are prime, we must have  $\gcd(\beta, q) = 1$ . We can apply Euler's theorem using modulus  $q$  to conclude  $\llbracket \beta^{\varphi(q)} \rrbracket_q = 1$ , so

$$\llbracket \beta^{1+y\varphi(n)} \rrbracket_q = \llbracket \beta \rrbracket_q \left( \llbracket \beta^{\varphi(q)} \rrbracket_q \right)^{y(p-1)} = \llbracket \beta \rrbracket_q.$$

We have two equations modulo  $p$  and modulo  $q$ , and can apply the Chinese Remainder Theorem (Thm. 10.5) to obtain one modulo  $pq$ . Hence

$$\llbracket \beta^{1+y\varphi(n)} \rrbracket_n = \llbracket \beta \rrbracket_n = \beta$$

as desired. In the case where  $\llbracket \beta \rrbracket_q = 0$ , the argument is the same after exchanging  $p$  and  $q$ .  $\square$

**Remarks.**

1. The symmetry between  $e$  and  $x$  means that we can encrypt with either the public or the private key, using the other for decryption. To use RSA works for authentication, the private key  $(n, x)$  is used to encrypt a message which can be decrypted with the public key  $(n, e)$ . Only one who possesses the private key could have encrypted such a message, which serves as a proof of identity.
2. If  $e$  is small, there is a possibility of having some message elements  $\beta$  for which  $\beta^e < n$ ; for such a  $\beta$ , decrypting without knowledge of  $\varphi(n)$  is easy: just compute the  $e$ -th root of  $\beta$ , since  $n$  is not a consideration. With care, one can either ensure that such  $\beta$  do not occur, or one can encrypt such messages in a slightly different manner.
3. For efficiency reasons, many implementations always use 3 for the encoding exponent  $e$ . In general, computing  $\llbracket \beta^e \rrbracket_n$  when  $e$  is of the form  $2^k + 1$  is quite efficient. Other common choices for  $e$  are 17 and 65537 ( $2^{16} + 1$ ), which are only slightly less efficient than 3.
4. Instead of using the Euler  $\varphi$ -function, one can use Carmichael's  $\lambda$ -function, which is the largest order of any element of  $\mathbb{Z}_n^*$ . In the case where  $n$  is a product of two primes  $p$  and  $q$ , we have  $\lambda(n) = \text{lcm}(p-1, q-1)$ . Since  $p-1$  and  $q-1$  are always even, the resulting private exponent  $x$  will be smaller.
5. There are a number of considerations in the choice of primes  $p$  and  $q$  that we will mostly ignore. For example,  $p$  and  $q$  should not be too close together, and both  $p-1$  and  $q-1$  should both have at least one large prime factor.

## 11.2 Implementing RSA in Maple

Now that we have covered how the RSA system works in theory, it is fairly straightforward to implement it.

### Implementing the basics of RSA

First, we write a Maple procedure which, given a size for  $n$  (our public base), will choose  $n$  as the product of two large primes  $p$  and  $q$ , as well as an exponent  $e$  which is relatively prime to  $\varphi(n)$ , and compute the decrypting exponent  $x$ .

For real security,  $n$  should be at least 200 decimal digits, meaning that the primes  $p$  and  $q$  will be about 100 digits each. Such large numbers are not easily factored (which is how we get our security), so how do we expect to be able to find 100-digit prime numbers? It turns out that there are a number of very efficient ways to test whether a number is prime without actually factoring it. Maple's `isprime` function uses such methods. So all we really need to do is generate a random number of around 100 digits, then use `nextprime` to get the next prime bigger than that number. We should use `randomize()` to ensure that we get different numbers each time (See the discussion of pseudo-random numbers in §5.2, page 4:17).

In `RSAsSetup` below, the approximate number of digits for  $n$  is given, and  $p$  is chosen as a prime with one less than half that many digits, while  $q$  has one more than half the number of digits of  $n$ . This means  $n$  must have at least 5 digits (or else  $p$  would have to be 2, 3, 5 or 7). After calculating  $n$  and  $\varphi(n)$ , a random number is chosen for the public exponent  $e$ . If this number is not relatively prime to  $\varphi(n)$ , another random choice is taken until such a number is found. Finally,  $x$  is computed as the inverse of  $e$  modulo  $\varphi(n)$ . The result of the function is the public key  $(n, e)$  and the private key  $(n, x)$ .

```
> randomize():

RSAsSetup:=proc(numdigits::posint)
  local p,q,phi,e,n, x, pm, qm;

  if (numdigits <= 4) then
    error("Too few digits for this to work");
  fi;

  pm:=floor((numdigits-1)/2);
  qm:=numdigits-pm;
  p:=nextprime(rand(10^(pm-1)..10^(pm))());
  q:=nextprime(rand(10^(qm-1)..10^(qm))());

  n:=p*q;
  phi:=(p-1)*(q-1);

  e:=rand(3..phi/2)();
  while(gcd(e,phi) <> 1) do
    e:=rand(3..phi/2)();
  od;
  x := modp(1/e, phi);
  return([[n,e],[n,x]]);
end:
```

We should remark that in certain cases, both  $p$  and  $q$  might be chosen at the extreme ends of their ranges, giving  $n$  with one fewer or one more digit than requested. If this were important, it could be repaired with a slightly more complicated program.

Once the keys are chosen, writing the heart of the RSA encoding is quite easy. We just need a function that eats a number  $\beta$  and a key  $(n, e)$ , and computes  $\llbracket \beta^e \rrbracket_n$ .

```
> RSAEncodeNum := proc(num::nonnegint, key::list(posint))
  return( modp( num &^ key[2], key[1] ));
end:
```

As a brief example, we'll generate a random pair of keys with a 5-digit modulus, and then encrypt and decrypt the number 47.

```
> keys:=RSAsSetup(5):
  public := keys[1];
  private:= keys[2];
```



```
public := [10793, 2371]
```

```
private := [10793, 31]
```

```
> RSAEncodeNum(47,public);
   RSAEncodeNum(% ,private);
```

```
3991
```

```
47
```

Of course, we want to encrypt more than a single number, but this is the heart of the whole thing.

### Making it Useful

In order to encrypt a string of text, we have to convert our text to a list of numbers. It is probably best to use  $k$ -graphs to represent our text, as in §5.3, where  $k$  is the largest integer power of the size of our alphabet that is smaller than  $n$ . While it will work to use a single character alphabet, this can compromise security significantly— see the remarks at the end of §11.1.

We can use `StringToKgraph` and `KgraphToString` from §5.3 to represent our strings as  $k$ -graphs. In order to determine the appropriate choice of  $k$ , we find the largest integer  $k$  so that

$$\text{Alen}^k < n$$

where `Alen` is the length of our alphabet.

Once we have our plaintext represented as a list of integers, we merely map the function `RSAEncodeNum` onto that list.

```
> Alphabet := cat("\n\t",Select(IsPrintable,convert([seq(i,i=1..255)],bytes))):
   Alen := length(Alphabet);
   k:=floor(log[Alen](public[1]));
```

```
Alen := 97
```

```
k := 2
```

```
> text := "Once upon a midnight dreary, while I pondered, weak and weary":
```

```

> StringToKgraph(text,k);

[7809, 6956, 8441, 7939, 274, 261, 7354, 7830, 7156, 8416, 6792, 6971, 8215, 1449,
8635, 7349, 6965, 4173, 7956, 7841, 6957, 6971, 1428, 8635, 6570, 271, 7827,
264, 6976, 8215, 91]

> crypt:=map(RSAEncodeNum,%,public);

crypt := [23972, 14519, 16728, 21226, 15537, 16168, 23547, 21311, 657, 15151, 20048,
7796, 961, 8740, 4392, 5615, 18965, 18341, 18628, 14070, 9286, 7796, 23471,
4392, 17418, 12708, 3247, 11788, 23561, 961, 9606]

```

We can reverse the process, using the private key, to decrypt.

```

> KgraphToString(
    map(RSAEncodeNum,crypt,private),
    k);

```

“Once upon a midnight dreary, while I pondered, weak and weary”

It is important to note that if we want to represent the result of the encryption as a string instead of a list of numbers, we must use  $(k + 1)$ -graphs rather than  $k$ -graphs. This is because  $\text{Alen}^k < n$ , so there will often be  $\beta$  so that  $\lceil \beta^e \rceil_n > \text{Alen}^k$ . In the example above, note that  $97^2 = 9409$ , and there are several elements of the list `crypt` larger than 9409.

We can put it all together in two simple procedures.

```

> RSAEncodeString := proc(text::string, key::list(posint))
    local Alen, k;
    global Alphabet;

    Alen := length(Alphabet);
    k:=floor(log[Alen](key[1]));
    KgraphToString(
        map(RSAEncodeNum,
            StringToKgraph(text, k),
            key),
        k+1);
end:

RSADecodeString := proc(text::string, key::list(posint))
    local Alen, k;
    global Alphabet;

    Alen := length(Alphabet);
    k:=floor(log[Alen](key[1]));
    KgraphToString(
        map(RSAEncodeNum,
            StringToKgraph(text, k+1),

```

```

        key),
    k);
end:

```

Now we try it out, with a 150-digit choice of  $n$ . Note that this means we'll be using 75-graphs, since  $97^{75} < 10^{150} < 97^{76}$ .

```

> keys:=RSAsetup(150):
   public:=keys[1]; private:=keys[2];

   public := [582757988242592241527296499231639188385179508302715855017\
1444848474216092727784907675976501184435973118498314592404067\
38402555037014297170960267438521, 18499822186252355491629282387\
3069724195125984703150098794787826856490657875937058983979216\
6564841269310221397039646221986692531497073722912893849813]

   private := [582757988242592241527296499231639188385179508302715855017\
1444848474216092727784907675976501184435973118498314592404067\
38402555037014297170960267438521, 33116290696193906028353889077\
81095383522640902859475195789426689481950851913681223087714308\
97081124962779673548200136089054119501927242337355077933565]

> crypt:=RSAEncodeString(text,public);

   crypt := “?!2<*h&l9![</*(=X5'oMT\\*Mnc /T%j*r~\“_=?b00'UP_\
gQ#O;oyVQdx3r—. N=s)!oidl?+?”

> RSADecodeString(crypt,private);

```

“Once upon a midnight dreary, while I pondered, weak and weary”

## 12 RSA encoding a file

Now we will write a procedure that uses RSA to encrypt a file, rather than a string. Of course, we could just read in the string and use `RSASecureString`, but in some cases, it makes sense to use a different alphabet for the crypttext. For example, we might want to accept all possible ASCII codes 0-255 as plaintext input, but restrict our encrypted output to use a limited collection of characters (so that it can be printed or emailed without transcription problems, say). For example, in the example below, we use the character set of RFC 2045<sup>35</sup>, which specifies how to encode base 64 MIME attachments for email.

---

<sup>35</sup><http://www.faqs.org/rfcs/rfc2045.html>

Here is our procedure which reads the contents of a file, encrypts it with RSA, and writes out the result to another file.

```

> Alphabet:=cat(Select(IsAlphaNumeric,convert([seq(i,i=1..255)],bytes)),"+/");
> RSAEncodeFile:=proc(plainfile::string, cryptfile::string, key::list(posint))
  local k, pf, cf, chunk, chunklen, codenum, crypt, Alen;
  global Alphabet;

  Alen    := length(Alphabet);
  chunklen:= ceil(log[Alen](key[1]));
  k       := floor(log[256](key[1]));

  try
    pf      := fopen(plainfile, READ, BINARY);
    cf      := fopen(cryptfile, WRITE, TEXT);
    chunk   := readbytes(pf, k);
    while (chunk <> 0) do
      codenum := map(RSAEncodeNum,
                    convert( chunk, base, 256, 256^k),
                    key);
      crypt:=KgraphToString(codenum, chunklen);
      writeline(cf, crypt);
      chunk := readbytes(pf, k);
    od;

  finally
    fclose(cryptfile);
    fclose(plainfile);
  end try;
  RETURN();
end:

> RSADecodeFile:=proc(cryptfile::string, plainfile::string, key::list(posint))
  local k, pf, cf, chunk, chunklen, plain, crypt, Alen;
  global Alphabet;

  Alen    := length(Alphabet);
  chunklen:= ceil(log[Alen](key[1]));
  k:=      floor(log[256](key[1]));

  try
    pf      := fopen(plainfile, WRITE, BINARY);
    cf      := fopen(cryptfile, READ, TEXT);
    chunk   := readline(cf);
    while (chunk <> 0) do
      plain:= convert(
        map(RSAEncodeNum,
            StringToKgraph(chunk, chunklen),
            key),
        base, 256^k, 256);
      writebytes(pf, plain);
    od;
  end try;
end:

```

```
        chunk := readline(cf);
    od;

    finally
        fclose(cryptfile);
        fclose(plainfile);
    end try;

    RETURN();
end;
```

## References

- [DH76] W. Diffie and M. Hellman, “New Directions in Cryptography”, *IEEE Transactions on Information Theory* **IT-22** no.6 (1976), p. 644–654.
- [Hil1] L. Hill, “Cryptology in an Algebraic Alphabet”, *American Mathematical Monthly* **36** (1929), p. 306–320.
- [Hil2] L. Hill, “Concerning the Linear Transformation Apparatus in Cryptography”, *American Mathematical Monthly* **38** (1931), p. 135–154.
- [HP] J. H. Humphreys & M. Y. Prest, *Numbers, Groups, and Codes*. Cambridge University Press, Cambridge. 1989.
- [Kahn] D. Kahn, *The Codebreakers; The Comprehensive History of Secret Communication from Ancient Times to the Internet*. Charles Scribner’s Sons, New York. 1996.
- [Kob] N. Koblitz, *A Course in Number Theory and Cryptography* (Graduate Texts in Mathematics, No 114). Springer-Verlag, New York. 1994.
- [NZM] I. Niven, H. Zuckerman, H. Montgomery, *An Introduction to the Theory of Numbers*, John Wiley & Sons, New York. 1991.
- [RSA] R. L. Rivest, A. Shamir and L. Adelman, “A method for obtaining digital signatures and public-key cryptosystems”, *Communications of the ACM* **21** (1978), p. 120–126.
- [Ros] K. Rosen, *Elementary Number Theory and its Applications*, 4th edition. Addison-Wesley, Boston. 2000.
- [Sch] B. Schneier, *Applied Cryptography: Protocols, Algorithms, and Source Code in C*, 2nd edition. John Wiley & Sons, New York. 1995.

