

First-order differential equations **1**

1.1 Introduction

This book is a study of differential equations and their applications. A differential equation is a relationship between a function of time and its derivatives. The equations

$$\frac{dy}{dt} = 3y^2 \sin(t + y) \quad (\text{i})$$

and

$$\frac{d^3y}{dt^3} = e^{-y} + t + \frac{d^2y}{dt^2} \quad (\text{ii})$$

are both examples of differential equations. The order of a differential equation is the order of the highest derivative of the function y that appears in the equation. Thus (i) is a first-order differential equation and (ii) is a third-order differential equation. By a solution of a differential equation we will mean a continuous function $y(t)$ which together with its derivatives satisfies the relationship. For example, the function

$$y(t) = 2 \sin t - \frac{1}{3} \cos 2t$$

is a solution of the second-order differential equation

$$\frac{d^2y}{dt^2} + y = \cos 2t$$

since

$$\begin{aligned} \frac{d^2}{dt^2} (2 \sin t - \frac{1}{3} \cos 2t) + (2 \sin t - \frac{1}{3} \cos 2t) \\ = (-2 \sin t + \frac{4}{3} \cos 2t) + 2 \sin t - \frac{1}{3} \cos 2t = \cos 2t. \end{aligned}$$

Differential equations appear naturally in many areas of science and the humanities. In this book, we will present serious discussions of the applications of differential equations to such diverse and fascinating problems as the detection of art forgeries, the diagnosis of diabetes, the increase in the percentage of sharks present in the Mediterranean Sea during World War I, and the spread of gonorrhoea. Our purpose is to show how researchers have used differential equations to solve, or try to solve, *real life* problems. And while we will discuss some of the great success stories of differential equations, we will also point out their limitations and document some of their failures.

1.2 First-order linear differential equations

We begin by studying first-order differential equations and we will assume that our equation is, or can be put, in the form

$$\frac{dy}{dt} = f(t, y). \quad (1)$$

The problem before us is this: Given $f(t, y)$ find all functions $y(t)$ which satisfy the differential equation (1). We approach this problem in the following manner. A fundamental principle of mathematics is that the way to solve a new problem is to reduce it, in some manner, to a problem that we have already solved. In practice this usually entails successively simplifying the problem until it resembles one we have already solved. Since we are presently in the business of solving differential equations, it is advisable for us to take inventory and list all the differential equations we can solve. If we assume that our mathematical background consists of just elementary calculus then the very sad fact is that the only first-order differential equation we can solve at present is

$$\frac{dy}{dt} = g(t) \quad (2)$$

where g is any integrable function of time. To solve Equation (2) simply integrate both sides with respect to t , which yields

$$y(t) = \int g(t) dt + c.$$

Here c is an arbitrary constant of integration, and by $\int g(t) dt$ we mean an anti-derivative of g , that is, a function whose derivative is g . Thus, to solve any other differential equation we must somehow reduce it to the form (2). As we will see in Section 1.9, this is impossible to do in most cases. Hence, we will not be able, without the aid of a computer, to solve most differential equations. It stands to reason, therefore, that to find those differential equations that we *can* solve, we should start with very simple equations

and not ones like

$$\frac{dy}{dt} = e^{\sin(t-37\sqrt{|y|})}$$

(which incidentally, cannot be solved exactly). Experience has taught us that the “simplest” equations are those which are *linear* in the dependent variable y .

Definition. The general first-order linear differential equation is

$$\frac{dy}{dt} + a(t)y = b(t). \quad (3)$$

Unless otherwise stated, the functions $a(t)$ and $b(t)$ are assumed to be continuous functions of time. We single out this equation and call it linear because the dependent variable y appears by itself, that is, no terms such as e^{-y} , y^3 or $\sin y$ etc. appear in the equation. For example $dy/dt = y^2 + \sin t$ and $dy/dt = \cos y + t$ are both *nonlinear* equations because of the y^2 and $\cos y$ terms respectively.

Now it is not immediately apparent how to solve Equation (3). Thus, we simplify it even further by setting $b(t) = 0$.

Definition. The equation

$$\frac{dy}{dt} + a(t)y = 0 \quad (4)$$

is called the *homogeneous* first-order linear differential equation, and Equation (3) is called the *nonhomogeneous* first-order linear differential equation for $b(t)$ not identically zero.

Fortunately, the homogeneous equation (4) can be solved quite easily. First, divide both sides of the equation by y and rewrite it in the form

$$\frac{\frac{dy}{dt}}{y} = -a(t).$$

Second, observe that

$$\frac{\frac{dy}{dt}}{y} \equiv \frac{d}{dt} \ln|y(t)|$$

where by $\ln|y(t)|$ we mean the natural logarithm of $|y(t)|$. Hence Equation (4) can be written in the form

$$\frac{d}{dt} \ln|y(t)| = -a(t). \quad (5)$$

But this is Equation (2) “essentially” since we can integrate both sides of (5) to obtain that

$$\ln|y(t)| = - \int a(t) dt + c_1$$

where c_1 is an arbitrary constant of integration. Taking exponentials of both sides yields

$$|y(t)| = \exp\left(- \int a(t) dt + c_1\right) = c \exp\left(- \int a(t) dt\right)$$

or

$$\left|y(t) \exp\left(\int a(t) dt\right)\right| = c. \quad (6)$$

Now, $y(t) \exp\left(\int a(t) dt\right)$ is a continuous function of time and Equation (6) states that its absolute value is constant. But if the absolute value of a continuous function $g(t)$ is constant then g itself must be constant. To prove this observe that if g is not constant, then there exist two different times t_1 and t_2 for which $g(t_1) = c$ and $g(t_2) = -c$. By the intermediate value theorem of calculus g must achieve all values between $-c$ and $+c$ which is impossible if $|g(t)| = c$. Hence, we obtain the equation $y(t) \exp\left(\int a(t) dt\right) = c$ or

$$y(t) = c \exp\left(- \int a(t) dt\right). \quad (7)$$

Equation (7) is said to be the *general solution* of the homogeneous equation since every solution of (4) must be of this form. Observe that an arbitrary constant c appears in (7). This should not be too surprising. Indeed, we will always expect an arbitrary constant to appear in the general solution of any first-order differential equation. To wit, if we are given dy/dt and we want to recover $y(t)$, then we must perform an integration, and this, of necessity, yields an arbitrary constant. Observe also that Equation (4) has infinitely many solutions; for each value of c we obtain a distinct solution $y(t)$.

Example 1. Find the general solution of the equation $(dy/dt) + 2ty = 0$.

Solution. Here $a(t) = 2t$ so that $y(t) = c \exp\left(- \int 2t dt\right) = c e^{-t^2}$.

Example 2. Determine the behavior, as $t \rightarrow \infty$, of all solutions of the equation $(dy/dt) + ay = 0$, a constant.

Solution. The general solution is $y(t) = c \exp\left(- \int a dt\right) = c e^{-at}$. Hence if $a < 0$, all solutions, with the exception of $y = 0$, approach infinity, and if $a > 0$, all solutions approach zero as $t \rightarrow \infty$.

In applications, we are usually not interested in all solutions of (4). Rather, we are looking for the *specific* solution $y(t)$ which at some initial time t_0 has the value y_0 . Thus, we want to determine a function $y(t)$ such that

$$\frac{dy}{dt} + a(t)y = 0, \quad y(t_0) = y_0. \quad (8)$$

Equation (8) is referred to as an initial-value problem for the obvious reason that of the totality of all solutions of the differential equation, we are looking for the one solution which initially (at time t_0) has the value y_0 . To find this solution we integrate both sides of (5) between t_0 and t . Thus

$$\int_{t_0}^t \frac{d}{ds} \ln|y(s)| ds = - \int_{t_0}^t a(s) ds$$

and, therefore

$$\ln|y(t)| - \ln|y(t_0)| = \ln \left| \frac{y(t)}{y(t_0)} \right| = - \int_{t_0}^t a(s) ds.$$

Taking exponentials of both sides of this equation we obtain that

$$\left| \frac{y(t)}{y(t_0)} \right| = \exp \left(- \int_{t_0}^t a(s) ds \right)$$

or

$$\left| \frac{y(t)}{y(t_0)} \exp \left(\int_{t_0}^t a(s) ds \right) \right| = 1.$$

The function inside the absolute value sign is a continuous function of time. Thus, by the argument given previously, it is either identically $+1$ or identically -1 . To determine which one it is, evaluate it at the point t_0 ; since

$$\frac{y(t_0)}{y(t_0)} \exp \left(\int_{t_0}^{t_0} a(s) ds \right) = 1$$

we see that

$$\frac{y(t)}{y(t_0)} \exp \left(\int_{t_0}^t a(s) ds \right) = 1.$$

Hence

$$y(t) = y(t_0) \exp \left(- \int_{t_0}^t a(s) ds \right) = y_0 \exp \left(- \int_{t_0}^t a(s) ds \right).$$

1 First-order differential equations

Example 3. Find the solution of the initial-value problem

$$\frac{dy}{dt} + (\sin t)y = 0, \quad y(0) = \frac{3}{2}.$$

Solution. Here $a(t) = \sin t$ so that

$$y(t) = \frac{3}{2} \exp\left(-\int_0^t \sin s \, ds\right) = \frac{3}{2} e^{(\cos t) - 1}.$$

Example 4. Find the solution of the initial-value problem

$$\frac{dy}{dt} + e^t y = 0, \quad y(1) = 2.$$

Solution. Here $a(t) = e^t$ so that

$$y(t) = 2 \exp\left(-\int_1^t e^{s^2} \, ds\right).$$

Now, at first glance this problem would seem to present a very serious difficulty in that we cannot integrate the function e^{s^2} directly. However, this solution is equally as valid and equally as useful as the solution to Example 3. The reason for this is twofold. First, there are very simple numerical schemes to evaluate the above integral to any degree of accuracy with the aid of a computer. Second, even though the solution to Example 3 is given explicitly, we still cannot evaluate it at any time t without the aid of a table of trigonometric functions and some sort of calculating aid, such as a slide rule, electronic calculator or digital computer.

We return now to the nonhomogeneous equation

$$\frac{dy}{dt} + a(t)y = b(t).$$

It should be clear from our analysis of the homogeneous equation that the way to solve the nonhomogeneous equation is to express it in the form

$$\frac{d}{dt}(\text{“something”}) = b(t)$$

and then to integrate both sides to solve for “something”. However, the expression $(dy/dt) + a(t)y$ does not appear to be the derivative of some simple expression. The next logical step in our analysis therefore should be the following: Can we make the left hand side of the equation to be d/dt of “something”? More precisely, we can multiply both sides of (3) by any continuous function $\mu(t)$ to obtain the equivalent equation

$$\mu(t) \frac{dy}{dt} + a(t) \mu(t)y = \mu(t)b(t). \quad (9)$$

(By equivalent equations we mean that every solution of (9) is a solution of (3) and vice-versa.) Thus, can we *choose* $\mu(t)$ so that $\mu(t)(dy/dt) + a(t)\mu(t)y$ is the derivative of some simple expression? The answer to this question is yes, and is obtained by observing that

$$\frac{d}{dt} \mu(t)y = \mu(t) \frac{dy}{dt} + \frac{d\mu}{dt} y.$$

Hence, $\mu(t)(dy/dt) + a(t)\mu(t)y$ will be equal to the derivative of $\mu(t)y$ if and only if $d\mu(t)/dt = a(t)\mu(t)$. But this is a first-order linear homogeneous equation for $\mu(t)$, i.e. $(d\mu/dt) - a(t)\mu = 0$ which we already know how to solve, and since we only need one such function $\mu(t)$ we set the constant c in (7) equal to one and take

$$\mu(t) = \exp\left(\int a(t) dt\right).$$

For this $\mu(t)$, Equation (9) can be written as

$$\frac{d}{dt} \mu(t)y = \mu(t)b(t). \quad (10)$$

To obtain the general solution of the nonhomogeneous equation (3), that is, to find all solutions of the nonhomogeneous equation, we take the indefinite integral (anti-derivative) of both sides of (10) which yields

$$\mu(t)y = \int \mu(t)b(t) dt + c$$

or

$$y = \frac{1}{\mu(t)} \left(\int \mu(t)b(t) dt + c \right) = \exp\left(-\int a(t) dt\right) \left(\int \mu(t)b(t) dt + c \right). \quad (11)$$

Alternately, if we are interested in the specific solution of (3) satisfying the initial condition $y(t_0) = y_0$, that is, if we want to solve the initial-value problem

$$\frac{dy}{dt} + a(t)y = b(t), \quad y(t_0) = y_0$$

then we can take the definite integral of both sides of (10) between t_0 and t to obtain that

$$\mu(t)y - \mu(t_0)y_0 = \int_{t_0}^t \mu(s)b(s) ds$$

or

$$y = \frac{1}{\mu(t)} \left(\mu(t_0)y_0 + \int_{t_0}^t \mu(s)b(s) ds \right). \quad (12)$$

Remark 1. Notice how we used our knowledge of the solution of the homogeneous equation to find the function $\mu(t)$ which enables us to solve the nonhomogeneous equation. This is an excellent illustration of how we use our knowledge of the solution of a simpler problem to solve a harder problem.

Remark 2. The function $\mu(t) = \exp\left(\int a(t) dt\right)$ is called an *integrating factor* for the nonhomogeneous equation since after multiplying both sides by $\mu(t)$ we can immediately integrate the equation to find all solutions.

Remark 3. The reader should not memorize formulae (11) and (12). Rather, we will solve all nonhomogeneous equations by first multiplying both sides by $\mu(t)$, by writing the new left-hand side as the derivative of $\mu(t)y(t)$, and then by integrating both sides of the equation.

Remark 4. An alternative way of solving the initial-value problem $(dy/dt) + a(t)y = b(t)$, $y(t_0) = y_0$ is to find the general solution (11) of (3) and then use the initial condition $y(t_0) = y_0$ to evaluate the constant c . If the function $\mu(t)b(t)$ cannot be integrated directly, though, then we must take the definite integral of (10) to obtain (12), and this equation is then approximated numerically.

Example 5. Find the general solution of the equation $(dy/dt) - 2ty = t$.

Solution. Here $a(t) = -2t$ so that

$$\mu(t) = \exp\left(\int a(t) dt\right) = \exp\left(-\int 2t dt\right) = e^{-t^2}.$$

Multiplying both sides of the equation by $\mu(t)$ we obtain the equivalent equation

$$e^{-t^2}\left(\frac{dy}{dt} - 2ty\right) = te^{-t^2} \quad \text{or} \quad \frac{d}{dt}e^{-t^2}y = te^{-t^2}.$$

Hence,

$$e^{-t^2}y = \int te^{-t^2} dt + c = \frac{-e^{-t^2}}{2} + c$$

and

$$y(t) = -\frac{1}{2} + ce^{t^2}.$$

Example 6. Find the solution of the initial-value problem

$$\frac{dy}{dt} + 2ty = t, \quad y(1) = 2.$$

Solution. Here $a(t) = 2t$ so that

$$\mu(t) = \exp\left(\int a(t) dt\right) = \exp\left(\int 2t dt\right) = e^{t^2}.$$

Multiplying both sides of the equation by $\mu(t)$ we obtain that

$$e^{t^2}\left(\frac{dy}{dt} + 2ty\right) = te^{t^2} \quad \text{or} \quad \frac{d}{dt}(e^{t^2}y) = te^{t^2}.$$

Hence,

$$\int_1^t \frac{d}{ds} e^{s^2} y(s) ds = \int_1^t s e^{s^2} ds$$

so that

$$e^{s^2} y(s) \Big|_1^t = \frac{e^{s^2}}{2} \Big|_1^t.$$

Consequently,

$$e^{t^2} y - 2e = \frac{e^{t^2}}{2} - \frac{e}{2}$$

and

$$y = \frac{1}{2} + \frac{3e}{2} e^{-t^2} = \frac{1}{2} + \frac{3}{2} e^{1-t^2}.$$

Example 7. Find the solution of the initial-value problem

$$\frac{dy}{dt} + y = \frac{1}{1+t^2}, \quad y(2) = 3.$$

Solution. Here $a(t) = 1$, so that

$$\mu(t) = \exp\left(\int a(t) dt\right) = \exp\left(\int 1 dt\right) = e^t.$$

Multiplying both sides of the equation by $\mu(t)$ we obtain that

$$e^t \left(\frac{dy}{dt} + y \right) = \frac{e^t}{1+t^2} \quad \text{or} \quad \frac{d}{dt} e^t y = \frac{e^t}{1+t^2}.$$

Hence

$$\int_2^t \frac{d}{ds} e^s y(s) ds = \int_2^t \frac{e^s}{1+s^2} ds,$$

so that

$$e^t y - 3e^2 = \int_2^t \frac{e^s}{1+s^2} ds$$

and

$$y = e^{-t} \left[3e^2 + \int_2^t \frac{e^s}{1+s^2} ds \right].$$

EXERCISES

In each of Problems 1–7 find the general solution of the given differential equation.

1. $\frac{dy}{dt} + y \cos t = 0$

2. $\frac{dy}{dt} + y \sqrt{t} \sin t = 0$

1 First-order differential equations

3. $\frac{dy}{dt} + \frac{2t}{1+t^2}y = \frac{1}{1+t^2}$

4. $\frac{dy}{dt} + y = te^t$

5. $\frac{dy}{dt} + t^2y = 1$

6. $\frac{dy}{dt} + t^2y = t^2$

7. $\frac{dy}{dt} + \frac{t}{1+t^2}y = 1 - \frac{t^3}{1+t^4}y$

In each of Problems 8–14, find the solution of the given initial-value problem.

8. $\frac{dy}{dt} + \sqrt{1+t^2}y = 0, \quad y(0) = \sqrt{5}$

9. $\frac{dy}{dt} + \sqrt{1+t^2}e^{-t}y = 0, \quad y(0) = 1$

10. $\frac{dy}{dt} + \sqrt{1+t^2}e^{-t}y = 0, \quad y(0) = 0$

11. $\frac{dy}{dt} - 2ty = t, \quad y(0) = 1$

12. $\frac{dy}{dt} + ty = 1 + t, \quad y\left(\frac{3}{2}\right) = 0$

13. $\frac{dy}{dt} + y = \frac{1}{1+t^2}, \quad y(1) = 2$

14. $\frac{dy}{dt} - 2ty = 1, \quad y(0) = 1$

15. Find the general solution of the equation

$$(1+t^2)\frac{dy}{dt} + ty = (1+t^2)^{5/2}.$$

(Hint: Divide both sides of the equation by $1+t^2$.)

16. Find the solution of the initial-value problem

$$(1+t^2)\frac{dy}{dt} + 4ty = t, \quad y(1) = \frac{1}{4}.$$

17. Find a continuous solution of the initial-value problem

$$y' + y = g(t), \quad y(0) = 0$$

where

$$g(t) = \begin{cases} 2, & 0 \leq t \leq 1 \\ 0, & t > 1 \end{cases}.$$

18. Show that every solution of the equation $(dy/dt) + ay = be^{-ct}$ where a and c are positive constants and b is any real number approaches zero as t approaches infinity.

19. Given the differential equation $(dy/dt) + a(t)y = f(t)$ with $a(t)$ and $f(t)$ continuous for $-\infty < t < \infty$, $a(t) \geq c > 0$, and $\lim_{t \rightarrow \infty} f(t) = 0$, show that every solution tends to zero as t approaches infinity.

When we derived the solution of the nonhomogeneous equation we tacitly assumed that the functions $a(t)$ and $b(t)$ were continuous so that we could perform the necessary integrations. If either of these functions was discontinuous at a point t_1 , then we would expect that our solutions might be discontinuous at $t = t_1$. Problems 20–23 illustrate the variety of things that

may happen. In Problems 20–22 determine the behavior of all solutions of the given differential equation as $t \rightarrow 0$, and in Problem 23 determine the behavior of all solutions as $t \rightarrow \pi/2$.

$$20. \frac{dy}{dt} + \frac{1}{t}y = \frac{1}{t^2}$$

$$21. \frac{dy}{dt} + \frac{1}{\sqrt{t}}y = e^{\sqrt{t}/2}$$

$$22. \frac{dy}{dt} + \frac{1}{t}y = \cos t + \frac{\sin t}{t}$$

$$23. \frac{dy}{dt} + y \tan t = \sin t \cos t.$$

1.3 The Van Meegeren art forgeries

After the liberation of Belgium in World War II, the Dutch Field Security began its hunt for Nazi collaborators. They discovered, in the records of a firm which had sold numerous works of art to the Germans, the name of a banker who had acted as an intermediary in the sale to Goering of the painting “Woman Taken in Adultery” by the famed 17th century Dutch painter Jan Vermeer. The banker in turn revealed that he was acting on behalf of a third rate Dutch painter H. A. Van Meegeren, and on May 29, 1945 Van Meegeren was arrested on the charge of collaborating with the enemy. On July 12, 1945 Van Meegeren startled the world by announcing from his prison cell that he had never sold “Woman Taken in Adultery” to Goering. Moreover, he stated that this painting and the very famous and beautiful “Disciples at Emmaus”, as well as four other presumed Vermeers and two de Hooghs (a 17th century Dutch painter) were his own work. Many people, however, thought that Van Meegeren was only lying to save himself from the charge of treason. To prove his point, Van Meegeren began, while in prison, to forge the Vermeer painting “Jesus Amongst the Doctors” to demonstrate to the skeptics just how good a forger of Vermeer he was. The work was nearly completed when Van Meegeren learned that a charge of forgery had been substituted for that of collaboration. He, therefore, refused to finish and age the painting so that hopefully investigators would not uncover his secret of aging his forgeries. To settle the question an international panel of distinguished chemists, physicists and art historians was appointed to investigate the matter. The panel took x-rays of the paintings to determine whether other paintings were underneath them. In addition, they analyzed the pigments (coloring materials) used in the paint, and examined the paintings for certain signs of old age.

Now, Van Meegeren was well aware of these methods. To avoid detection, he scraped the paint from old paintings that were not worth much, just to get the canvas, and he tried to use pigments that Vermeer would have used. Van Meegeren also knew that old paint was extremely hard, and impossible to dissolve. Therefore, he very cleverly mixed a chemical, phenoformaldehyde, into the paint, and this hardened into bakelite when the finished painting was heated in an oven.

However, Van Meegeren was careless with several of his forgeries, and the panel of experts found traces of the modern pigment cobalt blue. In addition, they also detected the phenoformaldehyde, which was not discovered until the turn of the 19th century, in several of the paintings. On the basis of this evidence Van Meegeren was convicted, of forgery, on October 12, 1947 and sentenced to one year in prison. While in prison he suffered a heart attack and died on December 30, 1947.

However, even following the evidence gathered by the panel of experts, many people still refused to believe that the famed “Disciples at Emmaus” was forged by Van Meegeren. Their contention was based on the fact that the other alleged forgeries and Van Meegeren’s nearly completed “Jesus Amongst the Doctors” were of a very inferior quality. Surely, they said, the creator of the beautiful “Disciples at Emmaus” could not produce such inferior pictures. Indeed, the “Disciples at Emmaus” was certified as an authentic Vermeer by the noted art historian A. Bredius and was bought by the Rembrandt Society for \$170,000. The answer of the panel to these skeptics was that because Van Meegeren was keenly disappointed by his lack of status in the art world, he worked on the “Disciples at Emmaus” with the fierce determination of proving that he was better than a third rate painter. After producing such a masterpiece his determination was gone. Moreover, after seeing how easy it was to dispose of the “Disciples at Emmaus” he devoted less effort to his subsequent forgeries. This explanation failed to satisfy the skeptics. They demanded a thoroughly scientific and conclusive proof that the “Disciples at Emmaus” was indeed a forgery. This was done recently in 1967 by scientists at Carnegie Mellon University, and we would now like to describe their work.

The key to the dating of paintings and other materials such as rocks and fossils lies in the phenomenon of radioactivity discovered at the turn of the century. The physicist Rutherford and his colleagues showed that the atoms of certain “radioactive” elements are unstable and that within a given time period a fixed proportion of the atoms spontaneously disintegrate to form atoms of a new element. Because radioactivity is a property of the atom, Rutherford showed that the radioactivity of a substance is directly proportional to the number of atoms of the substance present. Thus, if $N(t)$ denotes the number of atoms present at time t , then dN/dt , the number of atoms that disintegrate per unit time is proportional to N , that is,

$$\frac{dN}{dt} = -\lambda N. \quad (1)$$

The constant λ which is positive, is known as the decay constant of the substance. The larger λ is, of course, the faster the substance decays. One measure of the rate of disintegration of a substance is its *half-life* which is defined as the time required for half of a given quantity of radioactive atoms to decay. To compute the half-life of a substance in terms of λ , assume that at time t_0 , $N(t_0) = N_0$. Then, the solution of the initial-value

problem $dN/dt = -\lambda N$, $N(t_0) = N_0$ is

$$N(t) = N_0 \exp\left(-\lambda \int_{t_0}^t ds\right) = N_0 e^{-\lambda(t-t_0)}$$

or $N/N_0 = \exp(-\lambda(t-t_0))$. Taking logarithms of both sides we obtain that

$$-\lambda(t-t_0) = \ln \frac{N}{N_0}. \quad (2)$$

Now, if $N/N_0 = \frac{1}{2}$ then $-\lambda(t-t_0) = \ln \frac{1}{2}$ so that

$$(t-t_0) = \frac{\ln 2}{\lambda} = \frac{0.6931}{\lambda}. \quad (3)$$

Thus, the half-life of a substance is $\ln 2$ divided by the decay constant λ . The dimension of λ , which we suppress for simplicity of writing, is reciprocal time. If t is measured in years then λ has the dimension of reciprocal years, and if t is measured in minutes, then λ has the dimension of reciprocal minutes. The half-lives of many substances have been determined and recorded. For example, the half-life of carbon-14 is 5568 years and the half-life of uranium-238 is 4.5 billion years.

Now the basis of “radioactive dating” is essentially the following. From Equation (2) we can solve for $t-t_0 = 1/\lambda \ln(N_0/N)$. If t_0 is the time the substance was initially formed or manufactured, then the age of the substance is $1/\lambda \ln(N_0/N)$. The decay constant λ is known or can be computed, in most instances. Moreover, we can usually evaluate N quite easily. Thus, if we knew N_0 we could determine the age of the substance. But this is the real difficulty of course, since we usually do not know N_0 . In some instances though, we can either determine N_0 indirectly, or else determine certain suitable ranges for N_0 , and such is the case for the forgeries of Van Meegeren.

We begin with the following well-known facts of elementary chemistry. Almost all rocks in the earth’s crust contain a small quantity of uranium. The uranium in the rock decays to another radioactive element, and that one decays to another and another, and so forth (see Figure 1) in a series of elements that results in lead, which is not radioactive. The uranium (whose half-life is over four billion years) keeps feeding the elements following it in the series, so that as fast as they decay, they are replaced by the elements before them.

Now, all paintings contain a small amount of the radioactive element lead-210 (^{210}Pb), and an even smaller amount of radium-226 (^{226}Ra), since these elements are contained in white lead (lead oxide), which is a pigment that artists have used for over 2000 years. For the analysis which follows, it is important to note that white lead is made from lead metal, which, in turn, is extracted from a rock called lead ore, in a process called smelting. In this process, the lead-210 in the ore goes along with the lead metal. However, 90–95% of the radium and its descendants are removed with

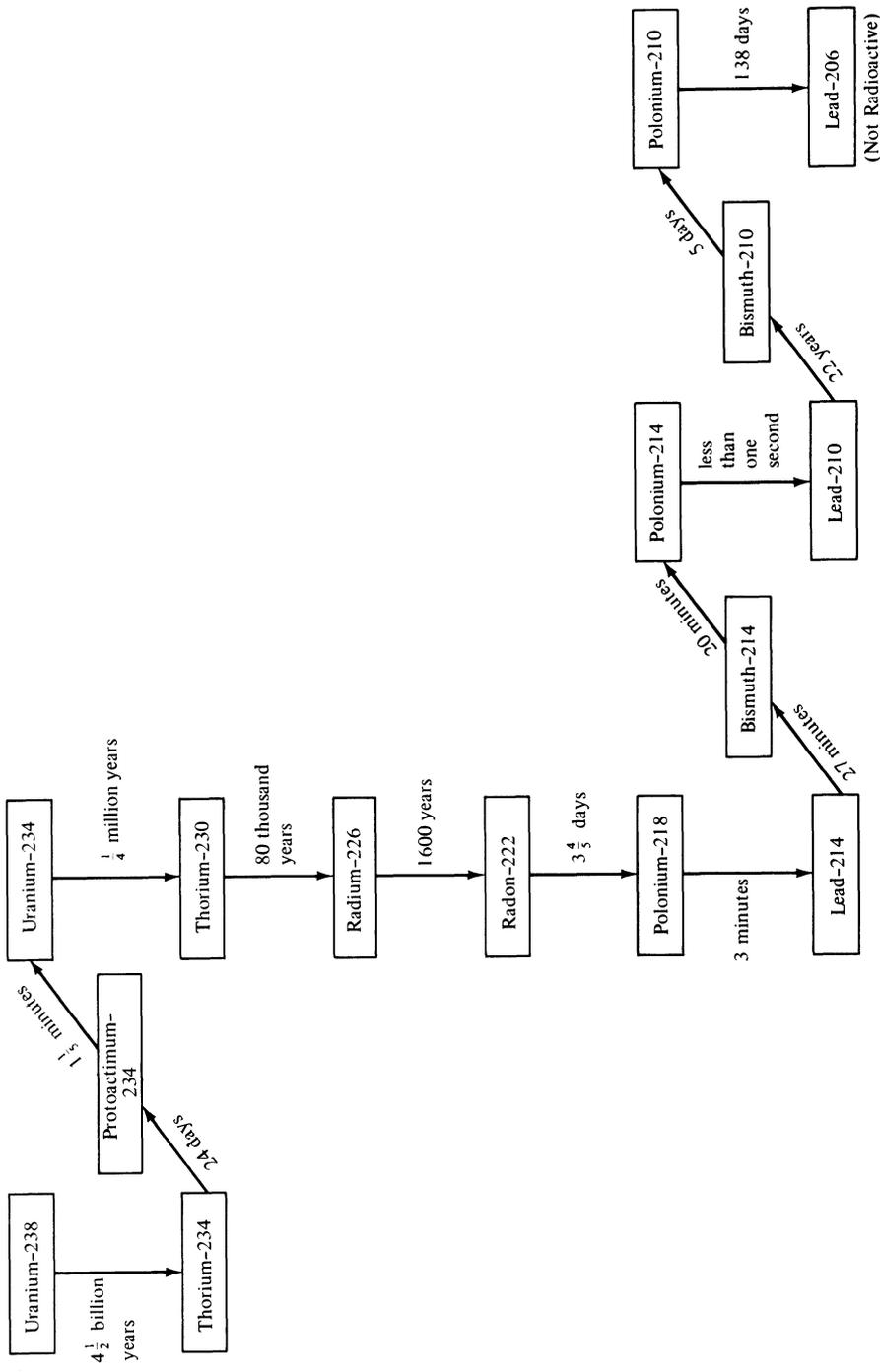


Figure 1. The Uranium series. (The times shown on the arrows are the half-lives of each step.)

other waste products in a material called slag. Thus, most of the supply of lead-210 is cut off and it begins to decay very rapidly, with a half-life of 22 years. This process continues until the lead-210 in the white lead is once more in radioactive equilibrium with the small amount of radium present, i.e. the disintegration of the lead-210 is exactly balanced by the disintegration of the radium.

Let us now use this information to compute the amount of lead-210 present in a sample in terms of the amount originally present at the time of manufacture. Let $y(t)$ be the amount of lead-210 per gram of white lead at time t , y_0 the amount of lead-210 per gram of white lead present at the time of manufacture t_0 , and $r(t)$ the number of disintegrations of radium-226 per minute per gram of white lead, at time t . If λ is the decay constant for lead-210, then

$$\frac{dy}{dt} = -\lambda y + r(t), \quad y(t_0) = y_0. \quad (4)$$

Since we are only interested in a time period of at most 300 years we may assume that the radium-226, whose half-life is 1600 years, remains constant, so that $r(t)$ is a constant r . Multiplying both sides of the differential equation by the integrating factor $\mu(t) = e^{\lambda t}$ we obtain that

$$\frac{d}{dt} e^{\lambda t} y = r e^{\lambda t}.$$

Hence

$$e^{\lambda t} y(t) - e^{\lambda t_0} y_0 = \frac{r}{\lambda} (e^{\lambda t} - e^{\lambda t_0})$$

or

$$y(t) = \frac{r}{\lambda} (1 - e^{-\lambda(t-t_0)}) + y_0 e^{-\lambda(t-t_0)}. \quad (5)$$

Now $y(t)$ and r can be easily measured. Thus, if we knew y_0 we could use Equation (5) to compute $(t - t_0)$ and consequently, we could determine the age of the painting. As we pointed out, though, we cannot measure y_0 directly. One possible way out of this difficulty is to use the fact that the original quantity of lead-210 was in radioactive equilibrium with the larger amount of radium-226 in the ore from which the metal was extracted. Let us, therefore, take samples of different ores and count the number of disintegrations of the radium-226 in the ores. This was done for a variety of ores and the results are given in Table 1 below. These numbers vary from 0.18 to 140. Consequently, the number of disintegrations of the lead-210 per minute per gram of white lead at the time of manufacture will vary from 0.18 to 140. This implies that y_0 will also vary over a very large interval, since the number of disintegrations of lead-210 is proportional to the amount present. Thus, we cannot use Equation (5) to obtain an accurate, or even a crude estimate, of the age of a painting.

Table 1. Ore and ore concentrate samples. All disintegration rates are per gram of white lead.

Description and Source	Disintegrations per minute of ^{226}Ra
Ore concentrate (Oklahoma-Kansas)	4.5
Crushed raw ore (S.E. Missouri)	2.4
Ore concentrate (S.E. Missouri)	0.7
Ore concentrate (Idaho)	2.2
Ore concentrate (Idaho)	0.18
Ore concentrate (Washington)	140.0
Ore concentrate (British Columbia)	1.9
Ore concentrate (British Columbia)	0.4
Ore concentrate (Bolivia)	1.6
Ore concentrate (Australia)	1.1

However, we can still use Equation (5) to distinguish between a 17th century painting and a modern forgery. The basis for this statement is the simple observation that if the paint is very old compared to the 22 year half-life of lead, then the amount of radioactivity from the lead-210 in the paint will be nearly equal to the amount of radioactivity from the radium in the paint. On the other hand, if the painting is modern (approximately 20 years old, or so) then the amount of radioactivity from the lead-210 will be much greater than the amount of radioactivity from the radium.

We make this argument precise in the following manner. Let us assume that the painting in question is either very new or about 300 years old. Set $t - t_0 = 300$ years in (5). Then, after some simple algebra, we see that

$$\lambda y_0 = \lambda y(t)e^{300\lambda} - r(e^{300\lambda} - 1). \quad (6)$$

If the painting is indeed a modern forgery, then λy_0 will be absurdly large. To determine what is an absurdly high disintegration rate we observe (see Exercise 1) that if the lead-210 decayed originally (at the time of manufacture) at the rate of 100 disintegrations per minute per gram of white lead, then the ore from which it was extracted had a uranium content of approximately 0.014 per cent. This is a fairly high concentration of uranium since the average amount of uranium in rocks of the earth's crust is about 2.7 parts per million. On the other hand, there are some very rare ores in the Western Hemisphere whose uranium content is 2-3 per cent. To be on the safe side, we will say that a disintegration rate of lead-210 is certainly absurd if it exceeds 30,000 disintegrations per minute per gram of white lead.

To evaluate λy_0 , we must evaluate the present disintegration rate, $\lambda y(t)$, of the lead-210, the disintegration rate r of the radium-226, and $e^{300\lambda}$. Since the disintegration rate of polonium-210 (^{210}Po) equals that of lead-210 after several years, and since it is easier to measure the disintegration rate of polonium-210, we substitute these values for those of lead-210. To compute

$e^{300\lambda}$, we observe from (3) that $\lambda = (\ln 2/22)$. Hence

$$e^{300\lambda} = e^{(300/22)\ln 2} = 2^{(150/11)}.$$

The disintegration rates of polonium-210 and radium-226 were measured for the “Disciples at Emmaus” and various other alleged forgeries and are given in Table 2 below.

Table 2. Paintings of questioned authorship. All disintegration rates are per minute, per gram of white lead.

Description	^{210}Po disintegration	^{226}Ra disintegration
“Disciples at Emmaus”	8.5	0.8
“Washing of Feet”	12.6	0.26
“Woman Reading Music”	10.3	0.3
“Woman Playing Mandolin”	8.2	0.17
“Lace Maker”	1.5	1.4
“Laughing Girl”	5.2	6.0

If we now evaluate λy_0 from (6) for the white lead in the painting “Disciples at Emmaus” we obtain that

$$\begin{aligned}\lambda y_0 &= (8.5)2^{150/11} - 0.8(2^{150/11} - 1) \\ &= 98,050\end{aligned}$$

which is unacceptably large. Thus, this painting must be a modern forgery. By a similar analysis, (see Exercises 2–4) the paintings “Washing of Feet”, “Woman Reading Music” and “Woman Playing Mandolin” were indisputably shown to be faked Vermeers. On the other hand, the paintings “Lace Maker” and “Laughing Girl” cannot be recently forged Vermeers, as claimed by some experts, since for these two paintings, the polonium-210 is very nearly in radioactive equilibrium with the radium-226, and no such equilibrium has been observed in any samples from 19th or 20th century paintings.

References

- Coremans, P., *Van Meergeren’s Faked Vermeers and De Hooghs*, Meulenhoff, Amsterdam, 1949.
- Keisch, B., Feller, R. L., Levine, A. S., Edwards, P. R., *Dating and Authenticating Works of Art by Measurement of Natural Alpha Emitters*, *Science* (155), 1238–1241, March 1967.
- Keisch, B., *Dating Works of Art through Their Natural Radioactivity: Improvements and Applications*, *Science*, 160, 413–415, April 1968.

1 First-order differential equations

EXERCISES

1. In this exercise we show how to compute the concentration of uranium in an ore from the dpm/(g of Pb) of the lead-210 in the ore.
 - (a) The half-life of uranium-238 is 4.51×10^9 years. Since this half-life is so large, we may assume that the amount of uranium in the ore is constant over a period of two to three hundred years. Let $N(t)$ denote the number of atoms of ^{238}U per gram of ordinary lead in the ore at time t . Since the lead-210 is in radioactive equilibrium with the uranium-238 in the ore, we know that $dN/dt = -\lambda N = -100$ dpm/g of Pb at time t_0 . Show that there are 3.42×10^{17} atoms of uranium-238 per gram of ordinary lead in the ore at time t_0 . (Hint: 1 year = 525,600 minutes.)
 - (b) Using the fact that one mole of uranium-238 weighs 238 grams, and that there are 6.02×10^{23} atoms in a mole, show that the concentration of uranium in the ore is approximately 0.014 percent.

For each of the paintings 2, 3, and 4 use the data in Table 2 to compute the disintegrations per minute of the original amount of lead-210 per gram of white lead, and conclude that each of these paintings is a forged Vermeer.

2. "Washing of Feet"
3. "Woman Reading Music"
4. "Woman Playing Mandolin"
5. The following problem describes a very accurate derivation of the age of uranium.
 - (a) Let $N_{238}(t)$ and $N_{235}(t)$ denote the number of atoms of ^{238}U and ^{235}U at time t in a given sample of uranium, and let $t=0$ be the time this sample was created. By the radioactive decay law,

$$\frac{d}{dt} N_{238}(t) = \frac{-\ln 2}{(4.5)10^9} N_{238}(t),$$

$$\frac{d}{dt} N_{235}(t) = \frac{-\ln 2}{0.707(10)^9} N_{235}(t).$$

Solve these equations for $N_{238}(t)$ and $N_{235}(t)$ in terms of their original numbers $N_{238}(0)$ and $N_{235}(0)$.

- (b) In 1946 the ratio of $^{238}\text{U}/^{235}\text{U}$ in any sample was 137.8. Assuming that equal amounts of ^{238}U and ^{235}U appeared in any sample at the time of its creation, show that the age of uranium is 5.96×10^9 years. This figure is universally accepted as the age of uranium.
6. In a samarskite sample discovered recently, there was 3 grams of Thorium (^{232}Th). Thorium decays to lead-208 (^{208}Pb) through the reaction $^{232}\text{Th} \rightarrow ^{208}\text{Pb} + 6(4^1\text{He})$. It was determined that 0.0376 of a gram of lead-208 was produced by the disintegration of the original Thorium in the sample. Given that the

half-life of Thorium is 13.9 billion years, derive the age of this samarskite sample. (Hint: 0.0376 grams of ^{208}Pb is the product of the decay of $(232/208) \times 0.0376$ grams of Thorium.)

One of the most accurate ways of dating archaeological finds is the method of carbon-14 (^{14}C) dating discovered by Willard Libby around 1949. The basis of this method is delightfully simple: The atmosphere of the earth is continuously bombarded by cosmic rays. These cosmic rays produce neutrons in the earth's atmosphere, and these neutrons combine with nitrogen to produce ^{14}C , which is usually called radiocarbon, since it decays radioactively. Now, this radiocarbon is incorporated in carbon dioxide and thus moves through the atmosphere to be absorbed by plants. Animals, in turn, build radiocarbon into their tissues by eating the plants. In living tissue, the rate of ingestion of ^{14}C exactly balances the rate of disintegration of ^{14}C . When an organism dies, though, it ceases to ingest carbon-14 and thus its ^{14}C concentration begins to decrease through disintegration of the ^{14}C present. Now, it is a fundamental assumption of physics that the rate of bombardment of the earth's atmosphere by cosmic rays has always been constant. This implies that the original rate of disintegration of the ^{14}C in a sample such as charcoal is the same as the rate measured today.* This assumption enables us to determine the age of a sample of charcoal. Let $N(t)$ denote the amount of carbon-14 present in a sample at time t , and N_0 the amount present at time $t=0$ when the sample was formed. If λ denotes the decay constant of ^{14}C (the half-life of carbon-14 is 5568 years) then $dN(t)/dt = -\lambda N(t)$, $N(0) = N_0$. Consequently, $N(t) = N_0 e^{-\lambda t}$. Now the present rate $R(t)$ of disintegration of the ^{14}C in the sample is given by $R(t) = \lambda N(t) = \lambda N_0 e^{-\lambda t}$ and the original rate of disintegration is $R(0) = \lambda N_0$. Thus, $R(t)/R(0) = e^{-\lambda t}$ so that $t = (1/\lambda) \ln[R(0)/R(t)]$. Hence if we measure $R(t)$, the present rate of disintegration of the ^{14}C in the charcoal, and observe that $R(0)$ must equal the rate of disintegration of the ^{14}C in a comparable amount of living wood, then we can compute the age t of the charcoal. The following two problems are real life illustrations of this method.

7. Charcoal from the occupation level of the famous Lascaux Cave in France gave an average count in 1950 of 0.97 disintegrations per minute per gram. Living wood gave 6.68 disintegrations. Estimate the date of occupation and hence the probable date of the remarkable paintings in the Lascaux Cave.
8. In the 1950 excavation at Nippur, a city of Babylonia, charcoal from a roof beam gave a count of 4.09 disintegrations per minute per gram. Living wood gave 6.68 disintegrations. Assuming that this charcoal was formed during the time of Hammurabi's reign, find an estimate for the likely time of Hamurabi's succession.

*Since the mid 1950's the testing of nuclear weapons has significantly increased the amount of radioactive carbon in our atmosphere. Ironically this unfortunate state of affairs provides us with yet another extremely powerful method of detecting art forgeries. To wit, many artists' materials, such as linseed oil and canvas paper, come from plants and animals, and so will contain the same concentration of carbon-14 as the atmosphere at the time the plant or animal dies. Thus linseed oil (which is derived from the flax plant) that was produced during the last few years will contain a much greater concentration of carbon-14 than linseed oil produced before 1950.

1.4 Separable equations

We solved the first-order linear homogeneous equation

$$\frac{dy}{dt} + a(t)y = 0 \quad (1)$$

by dividing both sides of the equation by $y(t)$ to obtain the equivalent equation

$$\frac{1}{y(t)} \frac{dy(t)}{dt} = -a(t) \quad (2)$$

and observing that Equation (2) can be written in the form

$$\frac{d}{dt} \ln|y(t)| = -a(t). \quad (3)$$

We then found $\ln|y(t)|$, and consequently $y(t)$, by integrating both sides of (3). In an exactly analogous manner, we can solve the more general differential equation

$$\frac{dy}{dt} = \frac{g(t)}{f(y)} \quad (4)$$

where f and g are continuous functions of y and t . This equation, and any other equation which can be put into this form, is said to be separable. To solve (4), we first multiply both sides by $f(y)$ to obtain the equivalent equation

$$f(y) \frac{dy}{dt} = g(t). \quad (5)$$

Then, we observe that (5) can be written in the form

$$\frac{d}{dt} F(y(t)) = g(t) \quad (6)$$

where $F(y)$ is any anti-derivative of $f(y)$; i.e., $F(y) = \int f(y) dy$. Consequently,

$$F(y(t)) = \int g(t) dt + c \quad (7)$$

where c is an arbitrary constant of integration, and we solve for $y = y(t)$ from (7) to find the general solution of (4).

Example 1. Find the general solution of the equation $dy/dt = t^2/y^2$.

Solution. Multiplying both sides of this equation by y^2 gives

$$y^2 \frac{dy}{dt} = t^2, \quad \text{or} \quad \frac{d}{dt} \frac{y^3(t)}{3} = t^2.$$

Hence, $y^3(t) = t^3 + c$ where c is an arbitrary constant, and $y(t) = (t^3 + c)^{1/3}$.

Example 2. Find the general solution of the equation

$$e^y \frac{dy}{dt} - t - t^3 = 0.$$

Solution. This equation can be written in the form

$$\frac{d}{dt} e^{y(t)} = t + t^3$$

and thus $e^{y(t)} = t^2/2 + t^4/4 + c$. Taking logarithms of both sides of this equation gives $y(t) = \ln(t^2/2 + t^4/4 + c)$.

In addition to the differential equation (4), we will often impose an initial condition on $y(t)$ of the form $y(t_0) = y_0$. The differential equation (4) together with the initial condition $y(t_0) = y_0$ is called an initial-value problem. We can solve an initial-value problem two different ways. Either we use the initial condition $y(t_0) = y_0$ to solve for the constant c in (7), or else we integrate both sides of (6) between t_0 and t to obtain that

$$F(y(t)) - F(y_0) = \int_{t_0}^t g(s) ds. \quad (8)$$

If we now observe that

$$F(y) - F(y_0) = \int_{y_0}^y f(r) dr, \quad (9)$$

then we can rewrite (8) in the simpler form

$$\int_{y_0}^y f(r) dr = \int_{t_0}^t g(s) ds. \quad (10)$$

Example 3. Find the solution $y(t)$ of the initial-value problem

$$e^y \frac{dy}{dt} - (t + t^3) = 0, \quad y(1) = 1.$$

Solution. Method (i). From Example 2, we know that the general solution of this equation is $y = \ln(t^2/2 + t^4/4 + c)$. Setting $t = 1$ and $y = 1$ gives $1 = \ln(3/4 + c)$, or $c = e - 3/4$. Hence, $y(t) = \ln(e - 3/4 + t^2/2 + t^4/4)$.

Method (ii). From (10),

$$\int_1^y e^r dr = \int_1^t (s + s^3) ds.$$

Consequently,

$$e^y - e = \frac{t^2}{2} + \frac{t^4}{4} - \frac{1}{2} - \frac{1}{4}, \quad \text{and} \quad y(t) = \ln(e - 3/4 + t^2/2 + t^4/4).$$

Example 4. Solve the initial-value problem $dy/dt = 1 + y^2$, $y(0) = 0$.

Solution. Divide both sides of the differential equation by $1 + y^2$ to obtain

the equivalent equation $1/(1+y^2)dy/dt=1$. Then, from (10)

$$\int_0^y \frac{dr}{1+r^2} = \int_0^t ds.$$

Consequently, $\arctan y = t$, and $y = \tan t$.

The solution $y = \tan t$ of the above problem has the disturbing property that it goes to $\pm \infty$ at $t = \pm \pi/2$. And what's even more disturbing is the fact that there is nothing at all in this initial-value problem which even hints to us that there is any trouble at $t = \pm \pi/2$. The sad fact of life is that solutions of perfectly nice differential equations can go to infinity in finite time. Thus, solutions will usually exist only on a finite open interval $a < t < b$, rather than for all time. Moreover, as the following example shows, different solutions of the same differential equation usually go to infinity at different times.

Example 5. Solve the initial-value problem $dy/dt = 1 + y^2$, $y(0) = 1$.

Solution. From (10)

$$\int_1^y \frac{dr}{1+r^2} = \int_0^t ds.$$

Consequently, $\arctan y - \arctan 1 = t$, and $y = \tan(t + \pi/4)$. This solution exists on the open interval $-3\pi/4 < t < \pi/4$.

Example 6. Find the solution $y(t)$ of the initial-value problem

$$y \frac{dy}{dt} + (1 + y^2) \sin t = 0, \quad y(0) = 1.$$

Solution. Dividing both sides of the differential equation by $1 + y^2$ gives

$$\frac{y}{1+y^2} \frac{dy}{dt} = -\sin t.$$

Consequently,

$$\int_1^y \frac{r dr}{1+r^2} = \int_0^t -\sin s ds,$$

so that

$$\frac{1}{2} \ln(1+y^2) - \frac{1}{2} \ln 2 = \cos t - 1.$$

Solving this equation for $y(t)$ gives

$$y(t) = \pm (2e^{-4\sin^2 t/2} - 1)^{1/2}.$$

To determine whether we take the plus or minus branch of the square root, we note that $y(0)$ is positive. Hence,

$$y(t) = (2e^{-4\sin^2 t/2} - 1)^{1/2}$$

This solution is only defined when

$$2e^{-4\sin^2 t/2} \geq 1$$

or

$$e^{4\sin^2 t/2} \leq 2. \quad (11)$$

Since the logarithm function is monotonic increasing, we may take logarithms of both sides of (11) and still preserve the inequality. Thus, $4\sin^2 t/2 < \ln 2$, which implies that

$$\left| \frac{t}{2} \right| \leq \arcsin \frac{\sqrt{\ln 2}}{2}$$

Therefore, $y(t)$ only exists on the open interval $(-a, a)$ where

$$a = 2\arcsin[\sqrt{\ln 2}/2].$$

Now, this appears to be a new difficulty associated with nonlinear equations, since $y(t)$ just “disappears” at $t = \pm a$, without going to infinity. However, this apparent difficulty can be explained quite easily, and moreover, can even be anticipated, if we rewrite the differential equation above in the standard form

$$\frac{dy}{dt} = -\frac{(1+y^2)\sin t}{y}.$$

Notice that this differential equation is not defined when $y = 0$. Therefore, if a solution $y(t)$ achieves the value zero at some time $t = t^*$, then we cannot expect it to be defined for $t > t^*$. This is exactly what happens here, since $y(\pm a) = 0$.

Example 7. Solve the initial-value problem $dy/dt = (1+y)t$, $y(0) = -1$.

Solution. In this case, we cannot divide both sides of the differential equation by $1+y$, since $y(0) = -1$. However, it is easily seen that $y(t) = -1$ is one solution of this initial-value problem, and in Section 1.10 we show that it is the only solution. More generally, consider the initial-value problem $dy/dt = f(y)g(t)$, $y(t_0) = y_0$, where $f(y_0) = 0$. Certainly, $y(t) = y_0$ is one solution of this initial-value problem, and in Section 1.10 we show that it is the only solution if $\partial f/\partial y$ exists and is continuous.

Example 8. Solve the initial-value problem

$$(1 + e^y)dy/dt = \cos t, \quad y(\pi/2) = 3.$$

Solution. From (10),

$$\int_3^y (1 + e^r) dr = \int_{\pi/2}^t \cos s ds$$

so that $y + e^y = 2 + e^3 + \sin t$. This equation cannot be solved explicitly for y

1 First-order differential equations

as a function of t . Indeed, most separable equations cannot be solved explicitly for y as a function of t . Thus, when we say that

$$y + e^y = 2 + e^3 + \sin t$$

is the solution of this initial-value problem, we really mean that it is an implicit, rather than an explicit solution. This does not present us with any difficulties in applications, since we can always find $y(t)$ numerically with the aid of a digital computer (see Section 1.11).

Example 9. Find all solutions of the differential equation $dy/dt = -t/y$.

Solution. Multiplying both sides of the differential equation by y gives $y dy/dt = -t$. Hence

$$y^2 + t^2 = c^2. \quad (12)$$

Now, the curves (12) are *closed*, and we cannot solve for y as a *single-valued* function of t . The reason for this difficulty, of course, is that the differential equation is not defined when $y=0$. Nevertheless, the circles $t^2 + y^2 = c^2$ are perfectly well defined, even when $y=0$. Thus, we will call the circles $t^2 + y^2 = c^2$ *solution curves* of the differential equation

$$dy/dt = -t/y.$$

More generally, we will say that any curve defined by (7) is a solution curve of (4).

EXERCISES

In each of Problems 1–5, find the general solution of the given differential equation.

1. $(1+t^2)\frac{dy}{dt} = 1+y^2$. *Hint:* $\tan(x+y) = \frac{\tan x + \tan y}{1 - \tan x \tan y}$.

2. $\frac{dy}{dt} = (1+t)(1+y)$

3. $\frac{dy}{dt} = 1 - t + y^2 - ty^2$

4. $\frac{dy}{dt} = e^{t+y+3}$

5. $\cos y \sin t \frac{dy}{dt} = \sin y \cos t$

In each of Problems 6–12, solve the given initial-value problem, and determine the interval of existence of each solution.

6. $t^2(1+y^2) + 2y\frac{dy}{dt} = 0$, $y(0) = 1$

7. $\frac{dy}{dt} = \frac{2t}{y + yt^2}$, $y(2) = 3$

1 First-order differential equations

20. Consider the differential equation

$$\frac{dy}{dt} = \frac{t+y+1}{t-y+3}. \quad (*)$$

We could solve this equation if the constants 1 and 3 were not present. To eliminate these constants, we make the substitution $t = T + h$, $y = Y + k$.

(a) Determine h and k so that (*) can be written in the form $dY/dT = (T + Y)/(T - Y)$.

(b) Find the general solution of (*). (See Exercise 18).

21. (a) Prove that the differential equation

$$\frac{dy}{dt} = \frac{at + by + m}{ct + dy + n}$$

where a, b, c, d, m , and n are constants, can always be reduced to $dy/dt = (at + by)/(ct + dy)$ if $ad - bc \neq 0$.

(b) Solve the above equation in the special case that $ad = bc$.

Find the general solution of the following equations.

22. $(1 + t - 2y) + (4t - 3y - 6)dy/dt = 0$

23. $(t + 2y + 3) + (2t + 4y - 1)dy/dt = 0$

1.5 Population models

In this section we will study first-order differential equations which govern the growth of various species. At first glance it would seem impossible to model the growth of a species by a differential equation since the population of any species always changes by integer amounts. Hence the population of any species can never be a differentiable function of time. However, if a given population is very large and it is suddenly increased by one, then the change is very small compared to the given population. Thus, we make the approximation that large populations change continuously and even differentially with time.

Let $p(t)$ denote the population of a given species at time t and let $r(t, p)$ denote the difference between its birth rate and its death rate. If this population is isolated, that is, there is no net immigration or emigration, then dp/dt , the rate of change of the population, equals $rp(t)$. In the most simplistic model we assume that r is constant, that is, it does not change with either time or population. Then, we can write down the following differential equation governing population growth:

$$\frac{dp(t)}{dt} = ap(t), \quad a = \text{constant.}$$

This is a linear equation and is known as the Malthusian law of population growth. If the population of the given species is p_0 at time t_0 , then $p(t)$ satisfies the initial-value problem $dp(t)/dt = ap(t)$, $p(t_0) = p_0$. The solution of this initial-value problem is $p(t) = p_0 e^{a(t-t_0)}$. Hence any species satisfying the Malthusian law of population growth grows exponentially with time.

Now, we have just formulated a very simple model for population growth; so simple, in fact, that we have been able to solve it completely in a few lines. It is important, therefore, to see if this model, with its simplicity, has any relationship at all with reality. Let $p(t)$ denote the human population of the earth at time t . It was estimated that the earth's human population was increasing at an average rate of 2% per year during the period 1960–1970. Let us start in the middle of this decade on January 1, 1965, at which time the U.S. Department of Commerce estimated the earth's population to be 3.34 billion people. Then, $t_0 = 1965$, $p_0 = 3.34 \times 10^9$ and $a = .02$, so that

$$p(t) = (3.34)10^9 e^{.02(t-1965)}.$$

One way of checking the accuracy of this formula is to compute the time required for the population of the earth to double, and then compare it to the observed value of 35 years. Our formula predicts that the population of the earth will double every T years, where

$$e^{.02T} = 2.$$

Taking logarithms of both sides of this equation gives $.02T = \ln 2$, so that

$$T = 50 \ln 2 \approx 34.6 \text{ years.}$$

This is in excellent agreement with the observed value. On the other hand, though, let us look into the distant future. Our equation predicts that the earth's population will be 200,000 billion in the year 2515, 1,800,000 billion in the year 2625, and 3,600,000 billion in the year 2660. These are astronomical numbers whose significance is difficult to gauge. The total surface of this planet is approximately 1,860,000 billion square feet. Eighty percent of this surface is covered by water. Assuming that we are willing to live on boats as well as land, it is easy to see that by the year 2515 there will be only 9.3 square feet per person; by 2625 each person will have only one square foot on which to stand; and by 2660 we will be standing two deep on each other's shoulders.

It would seem therefore, that this model is unreasonable and should be thrown out. However, we cannot ignore the fact that it offers exceptional agreement in the past. Moreover, we have additional evidence that populations do grow exponentially. Consider the *Microtus Arvallis* Pall, a small rodent which reproduces very rapidly. We take the unit of time to be a month, and assume that the population is increasing at the rate of 40% per

1 First-order differential equations

month. If there are two rodents present initially at time $t=0$, then $p(t)$, the number of rodents at time t , satisfies the initial-value problem

$$dp(t)/dt = 0.4p, \quad p(0) = 2.$$

Consequently,

$$p(t) = 2e^{0.4t}. \quad (1)$$

Table 1 compares the observed population with the population calculated from Equation (1).

Table 1. The growth of *Microtus Arvallis* Pall.

Months	0	2	6	10
p Observed	2	5	20	109
p Calculated	2	4.5	22	109.1

As one can see, there is excellent agreement.

Remark. In the case of the *Microtus Arvallis* Pall, p observed is very accurate since the pregnancy period is three weeks and the time required for the census taking is considerably less. If the pregnancy period were very short then p observed could not be accurate since many of the pregnant rodents would have given birth before the census was completed.

The way out of our dilemma is to observe that linear models for population growth are satisfactory *as long as* the population is not too large. When the population gets extremely large though, these models cannot be very accurate, since they do not reflect the fact that individual members are now competing with each other for the limited living space, natural resources and food available. Thus, we must add a competition term to our linear differential equation. A suitable choice of a competition term is $-bp^2$, where b is a constant, since the statistical average of the number of encounters of two members per unit time is proportional to p^2 . We consider, therefore, the modified equation

$$\frac{dp}{dt} = ap - bp^2.$$

This equation is known as the logistic law of population growth and the numbers a, b are called the vital coefficients of the population. It was first introduced in 1837 by the Dutch mathematical-biologist Verhulst. Now, the constant b , in general, will be very small compared to a , so that if p is not too large then the term $-bp^2$ will be negligible compared to ap and the

population will grow exponentially. However, when p is very large, the term $-bp^2$ is no longer negligible, and thus serves to slow down the rapid rate of increase of the population. Needless to say, the more industrialized a nation is, the more living space it has, and the more food it has, the smaller the coefficient b is.

Let us now use the logistic equation to predict the future growth of an isolated population. If p_0 is the population at time t_0 , then $p(t)$, the population at time t , satisfies the initial-value problem

$$\frac{dp}{dt} = ap - bp^2, \quad p(t_0) = p_0.$$

This is a separable differential equation, and from Equation (10), Section 1.4,

$$\int_{p_0}^p \frac{dr}{ar - br^2} = \int_{t_0}^t ds = t - t_0.$$

To integrate the function $1/(ar - br^2)$ we resort to partial fractions. Let

$$\frac{1}{ar - br^2} \equiv \frac{1}{r(a - br)} = \frac{A}{r} + \frac{B}{a - br}.$$

To find A and B , observe that

$$\frac{A}{r} + \frac{B}{a - br} = \frac{A(a - br) + Br}{r(a - br)} = \frac{Aa + (B - bA)r}{r(a - br)}.$$

Therefore, $Aa + (B - bA)r = 1$. Since this equation is true for all values of r , we see that $Aa = 1$ and $B - bA = 0$. Consequently, $A = 1/a$, $B = b/a$, and

$$\begin{aligned} \int_{p_0}^p \frac{dr}{r(a - br)} &= \frac{1}{a} \int_{p_0}^p \left(\frac{1}{r} + \frac{b}{a - br} \right) dr \\ &= \frac{1}{a} \left[\ln \frac{p}{p_0} + \ln \left| \frac{a - bp_0}{a - bp} \right| \right] = \frac{1}{a} \ln \frac{p}{p_0} \left| \frac{a - bp_0}{a - bp} \right|. \end{aligned}$$

Thus,

$$a(t - t_0) = \ln \frac{p}{p_0} \left| \frac{a - bp_0}{a - bp} \right|. \quad (2)$$

Now, it is a simple matter to show (see Exercise 1) that

$$\frac{a - bp_0}{a - bp(t)}$$

is always positive. Hence,

$$a(t - t_0) = \ln \frac{p}{p_0} \frac{a - bp_0}{a - bp}.$$

Taking exponentials of both sides of this equation gives

$$e^{a(t-t_0)} = \frac{p}{p_0} \frac{a - bp_0}{a - bp},$$

or

$$p_0(a - bp)e^{a(t-t_0)} = (a - bp_0)p.$$

Bringing all terms involving p to the left-hand side of this equation, we see that

$$[a - bp_0 + bp_0e^{a(t-t_0)}]p(t) = ap_0e^{a(t-t_0)}.$$

Consequently,

$$p(t) = \frac{ap_0e^{a(t-t_0)}}{a - bp_0 + bp_0e^{a(t-t_0)}} = \frac{ap_0}{bp_0 + (a - bp_0)e^{-a(t-t_0)}}. \quad (3)$$

Let us now examine Equation (3) to see what kind of population it predicts. Observe that as $t \rightarrow \infty$,

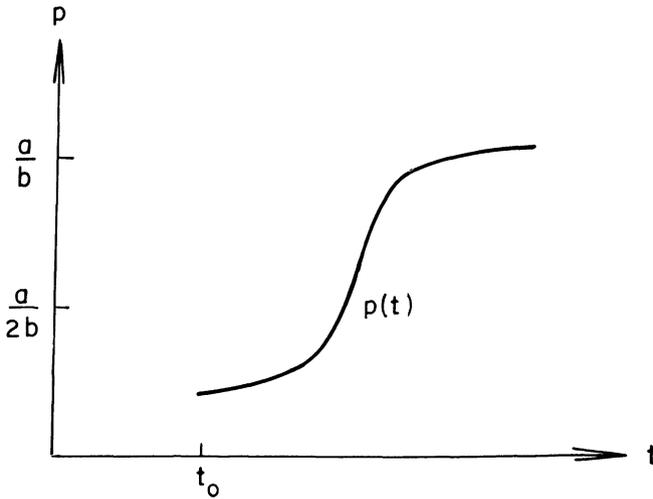
$$p(t) \rightarrow \frac{ap_0}{bp_0} = \frac{a}{b}.$$

Thus, *regardless of its initial value, the population always approaches the limiting value a/b* . Next, observe that $p(t)$ is a monotonically increasing function of time if $0 < p_0 < a/b$. Moreover, since

$$\frac{d^2p}{dt^2} = a \frac{dp}{dt} - 2bp \frac{dp}{dt} = (a - 2bp)p(a - bp),$$

we see that dp/dt is increasing if $p(t) < a/2b$, and that dp/dt is decreasing if $p(t) > a/2b$. Hence, if $p_0 < a/2b$, the graph of $p(t)$ must have the form given in Figure 1. Such a curve is called a logistic, or *S-shaped curve*. From its shape we conclude that the time period before the population reaches half its limiting value is a period of accelerated growth. After this point, the rate of growth decreases and in time reaches zero. This is a period of diminishing growth.

These predictions are borne out by an experiment on the protozoa *Paramecium caudatum* performed by the mathematical biologist G. F. Gause. Five individuals of *Paramecium* were placed in a small test tube containing 0.5 cm^3 of a nutritive medium, and for six days the number of individuals in every tube was counted daily. The *Paramecium* were found to increase at a rate of 230.9% per day when their numbers were low. The number of individuals increased rapidly at first, and then more slowly, until towards the fourth day it attained a maximum level of 375, saturating the test tube. From this data we conclude that if the *Paramecium caudatum* grow according to the logistic law $dp/dt = ap - bp^2$, then $a = 2.309$ and

Figure 1. Graph of $p(t)$

$b = 2.309/375$. Consequently, the logistic law predicts that

$$\begin{aligned}
 p(t) &= \frac{(2.309)5}{\frac{(2.309)5}{375} + \left(2.309 - \frac{(2.309)5}{375}\right)e^{-2.309t}} \\
 &= \frac{375}{1 + 74e^{-2.309t}}. \tag{4}
 \end{aligned}$$

(We have taken the initial time t_0 to be 0.) Figure 2 compares the graph of $p(t)$ predicted by Equation (4) with the actual measurements, which are denoted by o . As can be seen, the agreement is remarkably good.

In order to apply our results to predict the future human population of the earth, we must estimate the vital coefficients a and b in the logistic equation governing its growth. Some ecologists have estimated that the natural value of a is 0.029. We also know that the human population was increasing at the rate of 2% per year when the population was $(3.34)10^9$. Since $(1/p)(dp/dt) = a - bp$, we see that

$$0.02 = a - b(3.34)10^9.$$

Consequently, $b = 2.695 \times 10^{-12}$. Thus, according to the logistic law of population growth, the human population of the earth will tend to the limiting value of

$$\frac{a}{b} = \frac{0.029}{2.695 \times 10^{-12}} = 10.76 \text{ billion people}$$

Note that according to this prediction, we were still on the accelerated

1 First-order differential equations

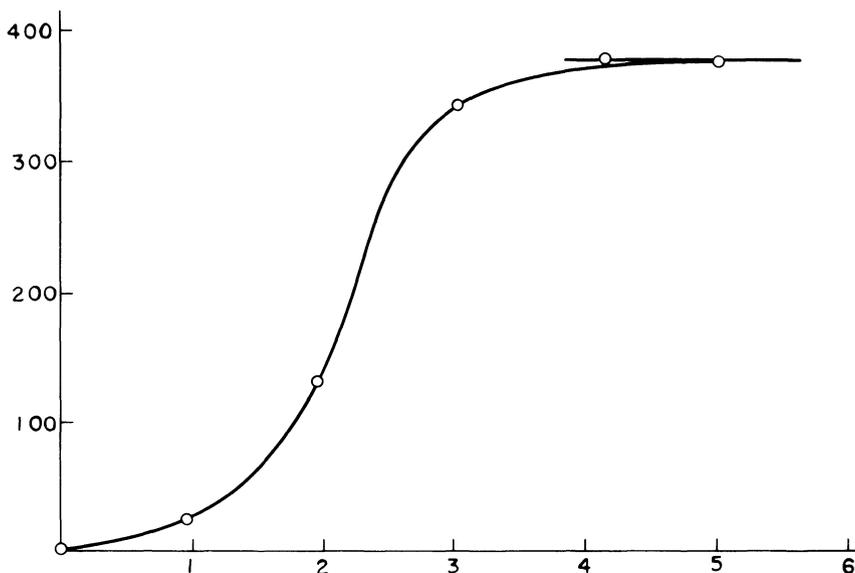


Figure 2. The growth of paramecium

growth portion of the logistic curve in 1965, since we had not yet attained half the limiting population predicted for us.

Remark. A student of mine once suggested that we use Equation (3) to find the time t at which $p(t) = 2$, and then we can deduce how long ago mankind appeared on earth. On the surface this seems like a fantastic idea. However, we cannot travel that far backwards into the past, since our model is no longer accurate when the population is small.

As another verification of the validity of the logistic law of population growth, we consider the equation

$$p(t) = \frac{197,273,000}{1 + e^{-0.03134(t - 1913.25)}} \quad (5)$$

which was introduced by Pearl and Reed as a model of the population growth of the United States. First, using the census figures for the years 1790, 1850, and 1910, Pearl and Reed found from (3) (see Exercise 2a) that $a = 0.03134$ and $b = (1.5887)10^{-10}$. Then (see Exercise 2b), Pearl and Reed calculated that the population of the United States reached half its limiting population of $a/b = 197,273,000$ in April 1913. Consequently (see Exercise 2c), we can rewrite (3) in the simpler form (5).

Table 2 below compares Pearl and Reed's predictions with the observed values of the population of the United States. These results are remarkable,

Table 2. Population of the U.S. from 1790–1950. (The last 4 entries were added by the Dartmouth College Writing Group.)

	Actual	Predicted	Error	%
1790	3,929,000	3,929,000	0	0.0
1800	5,308,000	5,336,000	28,000	0.5
1810	7,240,000	7,228,000	– 12,000	– 0.2
1820	9,638,000	9,757,000	119,000	1.2
1830	12,866,000	13,109,000	243,000	1.9
1840	17,069,000	17,506,000	437,000	2.6
1850	23,192,000	23,192,000	0	0.0
1860	31,443,000	30,412,000	– 1,031,000	– 3.3
1870	38,558,000	39,372,000	814,000	2.1
1880	50,156,000	50,177,000	21,000	0.0
1890	62,948,000	62,769,000	– 179,000	– 0.3
1900	75,995,000	76,870,000	875,000	1.2
1910	91,972,000	91,972,000	0	0.0
1920	105,711,000	107,559,000	1,848,000	1.7
1930	122,775,000	123,124,000	349,000	0.3
1940	131,669,000	136,653,000	4,984,000	3.8
1950	150,697,000	149,053,000	– 1,644,000	– 1.1

especially since we have not taken into account the large waves of immigration into the United States, and the fact that the United States was involved in five wars during this period.

In 1845 Verhulst prophesied a maximum population for Belgium of 6,600,000, and a maximum population for France of 40,000,000. Now, the population of Belgium in 1930 was already 8,092,000. This large discrepancy would seem to indicate that the logistic law of population growth is very inaccurate, at least as far as the population of Belgium is concerned. However, this discrepancy can be explained by the astonishing rise of industry in Belgium, and by the acquisition of the Congo which secured for the country sufficient additional wealth to support the extra population. Thus, after the acquisition of the Congo, and the astonishing rise of industry in Belgium, Verhulst should have lowered the vital coefficient b .

On the other hand, the population of France in 1930 was in remarkable agreement with Verhulst's forecast. Indeed, we can now answer the following tantalizing paradox: Why was the population of France increasing extremely slowly in 1930 while the French population of Canada was increasing very rapidly? After all, they are the same people! The answer to

this paradox, of course, is that the population of France in 1930 was very near its limiting value and thus was far into the period of diminishing growth, while the population of Canada in 1930 was still in the period of accelerated growth.

Remark 1. It is clear that technological developments, pollution considerations and sociological trends have significant influence on the vital coefficients a and b . Therefore, they must be re-evaluated every few years.

Remark 2. To derive more accurate models of population growth, we should not consider the population as made up of one homogeneous group of individuals. Rather, we should subdivide it into different age groups. We should also subdivide the population into males and females, since the reproduction rate in a population usually depends more on the number of females than on the number of males.

Remark 3. Perhaps the severest criticism leveled at the logistic law of population growth is that some populations have been observed to fluctuate periodically between two values, and any type of fluctuation is ruled out in a logistic curve. However, some of these fluctuations can be explained by the fact that when certain populations reach a sufficiently high density, they become susceptible to epidemics. The epidemic brings the population down to a lower value where it again begins to increase, until when it is large enough, the epidemic strikes again. In Exercise 10 we derive a model to describe this phenomenon, and we apply this model in Exercise 11 to explain the sudden appearance and disappearance of hordes of small rodents.

Epilog. The following article appeared in the New York Times on March 26, 1978, and was authored by Nick Eberstadt.

The gist of the following article is that it is very difficult, using statistical methods alone, to make accurate population projections even 30 years into the future. In 1970, demographers at the United Nations projected the population of the earth to be 6.5 billion people by the year 2000. Only six years later, this forecast was revised downward to 5.9 billion people.

Let us now use Equation (3) to predict the population of the earth in the year 2000. Setting $a = .029$, $b = 2.695 \times 10^{-12}$, $p_0 = 3.34 \times 10^9$, $t_0 = 1965$, and $t = 2,000$ gives

$$\begin{aligned} p(2000) &= \frac{(.029)(3.34)10^9}{.009 + (.02)e^{-(.029)35}} \\ &= \frac{29(3.34)}{9 + 20e^{-1.015}} 10^9 \\ &= 5.96 \text{ billion people!} \end{aligned}$$

This is another spectacular application of the logistic equation.

World Population Figures Are Misleading

The rate of world population growth has risen fairly steadily for most of man's history, but within the last decade it has peaked, and now appears to be declining. How has this happened, and why?

The "how" is fairly easy. It is not because famines and ecological catastrophes have elevated the death rates. Rather, a large and generally unexpected decrease in fertility in the less developed countries has taken place. From 1970 to 1977 birth rates in the less developed world (excluding China) fell from about 42 to nearly 36 per thousand. This is still higher than the 17 per thousand in developed countries, but the rate of fertility decline appears to be accelerating: the six-point drop of the past seven years compares with a two-point decline for the previous twenty.

This fertility decline has been a very uneven process. The average birth-rate drop of about 13 percent since 1970 in poor world birth rates reflects a very rapid decline in certain countries, while a great many others have remained almost totally unaffected.

Why fertility has dropped so rapidly in the past decade—and why it has dropped so dramatically in some places, but not in others—is far more difficult to explain. Demographers and sociologists offer explanations having to do with social change in the poor world. Unfortunately, these partial explanations are more often theory than tested

fact, and there seems to be an exception for almost every rule.

The debate over family planning is characteristic. Surveys show that in some nations as many as a fifth of the children were "mistakes" who presumably would not have been born if parents had had better contraceptives. Family planning experts such as Parker Mauldin of the Population Council have pointed out that no poor nation without an active family planning program has significantly lowered its fertility. On the other hand, such sociologists as William Petersen of Ohio State University attribute the population decline in these nations to social and economic development rather than increased contraceptive use, arguing that international "population control" programs have usually been clumsy and insensitive (or worse), and that in any event even a well-received change in contraceptive "technology" does not necessarily influence parents to want fewer children.

The effects of income distribution are less vociferously debated, but are almost as mysterious. James Kocher and Robert Repetto, both of Harvard, have argued that more equitable income distribution in less developed countries contributes to fertility decline. They have pointed out that such countries as Sri Lanka, South Korea, Cuba and China have seen their fertility rates fall as

their income distribution became more nearly equal. Improving a nation's income distribution, however, appears to be neither a necessary nor a sufficient condition for inducing a fertility drop. Income distribution in Burma, for example, has presumably equalized somewhat under 30 years of homemade socialism, but birth rates have hardly fallen at all, while Mexico and Colombia, with highly unequal income distributions, have found their birth rates plummeting in recent years.

One key to changes in fertility levels may be the economic costs and benefits from children. In peasant societies, where children are afforded few amenities and start work young, they may become economic assets early: A recent study in Bangladesh by Mead Cain of the Population Council put the age for boys at 12. Furthermore, children (or more precisely, sons) may also serve as social security and unemployment insurance when parents become too old and weak to work, or when work is unavailable. Where birth rates in poor countries are dropping, social and economic development may be making children less necessary as sources of income and security, but so little work has been done in this area that this is still just a reasonable speculation.

Some of the many other factors whose effects on fertility have been studied are urbanization, education, occupational struc-

ture, public health and the status of women. One area which population experts seem to have shied away from, however, is the non-quantifiable realm of attitudes, beliefs and values which may have had much to do with the recent changes in the decisions of hundreds of millions of couples. Cultural differences, ethnic conflicts, psychological, ideological and even political changes could clearly have effects on fertility. As Maris Vonovskis of the House Select Committee on Population has said, just because you can't measure something doesn't mean it isn't important.

What does the decline in fertility mean about future levels of population? Obviously, if the drop continues, population growth will be slower than previously anticipated, and world population will eventually stabilize at a lower level. Only five years ago the United Nations "medium variant" projection for world population in the year 2000 was 6.5 billion; last year this was dropped more than 200 million, and recent work by Gary Littman and Nathan Keyfitz at the Harvard Center for Population Studies shows that in the light of recent changes, one might easily drop it 400 million more.

Population projections, however, are a very tricky business. To begin with, the figures for today's population, upon which

tomorrow's projections must be based, contain large margins of error. For example, estimates for China's population run from 750 million to over 950 million. By the account of John Durand of the University of Pennsylvania, the margins of error for world population add up to over 200 million; historian Fernand Braudel puts the margin of error at 10 percent, which, given the world's approximate present population, means about 400 million people.

Population projections inspire even less confidence than population estimates, for they hinge on predicting birth and death rates for the future. These can change rapidly and unexpectedly: two extreme examples are Sri Lanka's 34 percent drop in the death rate in just two years, and Japan's 50 percent drop in the birth rate in 10. "Medium variant" U.N. projections computed just 17 years before 1975 overestimate Russia's population by 10 to 20 million, and underestimate India's by 50 million. Even projections for the United States done in 1966 overestimate its population only nine years later by over 10 million. Enormous as that gap may sound, it seems quite small next to those of the 1930's estimates which extrapolated low Depression era birth rates into an American population peaking at 170 million in the late 1970's (the population now is over

220 million), and then declining!

Could birth rates in the less developed world, which now appear to be declining at an accelerating pace, suddenly stabilize, or even rise again? This could theoretically happen. Here are four of the many reasons: 1) The many countries where fertility has as yet been unaffected by the decline might simply continue to be unaffected far into the future. 2) Since sterility and infertility are widespread in many of the poorest and most disease-ridden areas of the world, improvements in health and nutrition there could raise birth rates. 3) The Gandhi regime's cold-hearted and arbitrary mass sterilization regimen may have hardened that sixth of the world against future family limitation messages. 4) If John Aird of the Department of Commerce and others are correct that China's techniques of political mobilization and social persuasion have induced many parents to have fewer children than they actually want, a relaxation of these rules for whatever reasons might make the birth rate of China's enormous population rise. One of the only long-term rules about population projections which has held up is that within their limits of accuracy (about five years in the future) they can tell nothing interesting, and when they start giving interesting results, they are no longer accurate.

References

1. Gause, G. F., *The Struggle for Existence*, Dover Publications, New York, 1964.
2. Pearl and Reed, *Proceedings of the National Academy of Sciences*, 1920, p. 275.

EXERCISES

1. Prove that $(a - bp_0)/(a - bp(t))$ is positive for $t_0 < t < \infty$. *Hint*: Use Equation (2) to show that $p(t)$ can never equal a/b if $p_0 \neq a/b$.
2. (a) Choose 3 times t_0 , t_1 , and t_2 , with $t_1 - t_0 = t_2 - t_1$. Show that (3) determines a and b uniquely in terms of t_0 , $p(t_0)$, t_1 , $p(t_1)$, t_2 , and $p(t_2)$.
 (b) Show that the period of accelerated growth for the United States ended in April, 1913.
 (c) Let a population $p(t)$ grow according to the logistic law (3), and let \bar{t} be the time at which half the limiting population is achieved. Show that

$$p(t) = \frac{a/b}{1 + e^{-a(t-\bar{t})}}.$$

3. In 1879 and 1881 a number of yearling bass were seined in New Jersey, taken across the continent in tanks by train, and planted in San Francisco Bay. A total of only 435 Striped Bass survived the rigors of these two trips. Yet, in 1899, the commercial net catch alone was 1,234,000 pounds. Since the growth of this population was so fast, it is reasonable to assume that it obeyed the Malthusian law $dp/dt = ap$. Assuming that the average weight of a bass fish is three pounds, and that in 1899 every tenth bass fish was caught, find a lower bound for a .
4. Suppose that a population doubles its original size in 100 years, and triples it in 200 years. Show that this population cannot satisfy the Malthusian law of population growth.
5. Assume that $p(t)$ satisfies the Malthusian law of population growth. Show that the increases in p in successive time intervals of equal duration form the terms of a geometric progression. This is the source of Thomas Malthus' famous dictum "Population when unchecked increases in a geometrical ratio. Subsistence increases only in an arithmetic ratio. A slight acquaintance with numbers will show the immensity of the first power in comparison of the second."
6. A population grows according to the logistic law, with a limiting population of 5×10^8 individuals. When the population is low it doubles every 40 minutes. What will the population be after two hours if initially it is (a) 10^8 , (b) 10^9 ?
7. A family of salmon living off the Alaskan Coast obeys the Malthusian law of population growth $dp(t)/dt = 0.003p(t)$, where t is measured in minutes. At time $t = 0$ a group of sharks establishes residence in these waters and begins attacking the salmon. The rate at which salmon are killed by the sharks is $0.001p^2(t)$, where $p(t)$ is the population of salmon at time t . Moreover, since an undesirable element has moved into their neighborhood, 0.002 salmon per minute leave the Alaskan waters.
 (a) Modify the Malthusian law of population growth to take these two factors into account.

1 First-order differential equations

- (b) Assume that at time $t=0$ there are one million salmon. Find the population $p(t)$. What happens as $t \rightarrow \infty$?
- (c) Show that the above model is really absurd. Hint: Show, according to this model, that the salmon population decreases from one million to about one thousand in one minute.

8. The population of New York City would satisfy the logistic law

$$\frac{dp}{dt} = \frac{1}{25} p - \frac{1}{(25)10^6} p^2,$$

where t is measured in years, if we neglected the high emigration and homicide rates.

- (a) Modify this equation to take into account the fact that 9,000 people per year move from the city, and 1,000 people per year are murdered.
- (b) Assume that the population of New York City was 8,000,000 in 1970. Find the population for all future time. What happens as $t \rightarrow \infty$?
9. An initial population of 50,000 inhabits a microcosm with a carrying capacity of 100,000. After five years, the population has increased to 60,000. Show that the natural growth rate a for this population is $(1/5)\ln 3/2$.
10. We can model a population which becomes susceptible to epidemics in the following manner. Assume that our population is originally governed by the logistic law

$$\frac{dp}{dt} = ap - bp^2 \tag{i}$$

and that an epidemic strikes as soon as p reaches a certain value Q , with Q less than the limiting population a/b . At this stage the vital coefficients become $A < a$, $B < b$, and Equation (i) is replaced by

$$\frac{dp}{dt} = Ap - Bp^2. \tag{ii}$$

Suppose that $Q > A/B$. The population then starts decreasing. A point is reached when the population falls below a certain value $q > A/B$. At this moment the epidemic ceases and the population again begins to grow following Equation (i), until the incidence of a fresh epidemic. In this way there are periodic fluctuations of p between q and Q . We now indicate how to calculate the period T of these fluctuations.

- (a) Show that the time T_1 taken by the first part of the cycle, when p increases from q to Q is given by

$$T_1 = \frac{1}{a} \ln \frac{Q(a - bq)}{q(a - bQ)}.$$

- (b) Show that the time T_2 taken by the second part of the cycle, when p decreases from Q to q is given by

$$T_2 = \frac{1}{A} \ln \frac{q(QB - A)}{Q(qB - A)}.$$

Thus, the time for the entire cycle is $T_1 + T_2$.

11. It has been observed that plagues appear in mice populations whenever the population becomes too large. Further, a local increase of density attracts predators in large numbers. These two factors will succeed in destroying 97-98% of a population of small rodents in two or three weeks, and the density then falls to a level at which the disease cannot spread. The population, reduced to 2% of its maximum, finds its refuges from the predators sufficient, and its food abundant. The population therefore begins to grow again until it reaches a level favorable to another wave of disease and predation. Now, the speed of reproduction in mice is so great that we may set $b=0$ in Equation (i) of Exercise 7. In the second part of the cycle, on the contrary, A is very small in comparison with B , and it may be neglected therefore in Equation (ii).

(a) Under these assumptions, show that

$$T_1 = \frac{1}{a} \ln \frac{Q}{q} \quad \text{and} \quad T_2 = \frac{Q-q}{qQB}.$$

(b) Assuming that T_1 is approximately four years, and Q/q is approximately fifty, show that a is approximately one. This value of a , incidentally, corresponds very well with the rate of multiplication of mice in natural circumstances.

12. There are many important classes of organisms whose birth rate is *not* proportional to the population size. Suppose, for example, that each member of the population requires a partner for reproduction, and that each member relies on chance encounters for meeting a mate. If the expected number of encounters is proportional to the product of the numbers of males and females, and if these are equally distributed in the population, then the number of encounters, and hence the birthrate too, is proportional to p^2 . The death rate is still proportional to p . Consequently, the population size $p(t)$ satisfies the differential equation

$$\frac{dp}{dt} = bp^2 - ap, \quad a, b > 0.$$

Show that $p(t)$ approaches 0 as $t \rightarrow \infty$ if $p_0 < a/b$. Thus, once the population size drops below the critical size a/b , the population tends to extinction. Thus, a species is classified endangered if its current size is perilously close to its critical size.

1.6 The spread of technological innovations

Economists and sociologists have long been concerned with how a technological change, or innovation, spreads in an industry. Once an innovation is introduced by one firm, how soon do others in the industry come to adopt it, and what factors determine how rapidly they follow? In this section we construct a model of the spread of innovations among farmers, and then show that this same model also describes the spread of innovations in such diverse industries as bituminous coal, iron and steel, brewing, and railroads.

Assume that a new innovation is introduced into a fixed community of N farmers at time $t=0$. Let $p(t)$ denote the number of farmers who have

adopted at time t . As in the previous section, we make the approximation that $p(t)$ is a continuous function of time, even though it obviously changes by integer amounts. The simplest realistic assumption that we can make concerning the spread of this innovation is that a farmer adopts the innovation only after he has been told of it by a farmer who has already adopted. Then, the number of farmers Δp who adopt the innovation in a small time interval Δt is directly proportional to the number of farmers p who have already adopted, and the number of farmers $N - p$ who are as yet unaware. Hence, $\Delta p = cp(N - p)\Delta t$ or $\Delta p/\Delta t = cp(N - p)$ for some positive constant c . Letting $\Delta t \rightarrow 0$, we obtain the differential equation

$$\frac{dp}{dt} = cp(N - p). \quad (1)$$

This is the logistic equation of the previous section if we set $a = cN$, $b = c$. Assuming that $p(0) = 1$; i.e., one farmer has adopted the innovation at time $t = 0$, we see that $p(t)$ satisfies the initial-value problem

$$\frac{dp}{dt} = cp(N - p), \quad p(0) = 1. \quad (2)$$

The solution of (2) is

$$p(t) = \frac{Ne^{cNt}}{N - 1 + e^{cNt}} \quad (3)$$

which is a logistic function (see Section 1.5). Hence, our model predicts that the adoption process accelerates up to that point at which half the community is aware of the innovation. After this point, the adoption process begins to decelerate until it eventually reaches zero.

Let us compare the predictions of our model with data on the spread of two innovations through American farming communities in the middle 1950's. Figure 1 represents the cumulative number of farmers in Iowa during 1944–1955 who adopted 2,4-D weed spray, and Figure 2 represents the cumulative percentage of corn acreage in hybrid corn in three American states during the years 1934–1958. The circles in these figures are the actual measurements, and the graphs were obtained by connecting these measurements with straight lines. As can be seen, these curves have all the properties of logistic curves, and on the whole, offer very good agreement with our model. However, there are two discrepancies. First, the actual point at which the adoption process ceases to accelerate is not always when fifty per cent of the population has adopted the innovation. As can be seen from Figure 2, the adoption process for hybrid corn began to decelerate in Alabama only after nearly sixty per cent of the farmers had adopted the innovation. Second, the agreement with our model is much better in the later stages of the adoption process than in the earlier stages.

The source of the second discrepancy is our assumption that a farmer only learns of an innovation through contact with another farmer. This is not entirely true. Studies have shown that mass communication media such

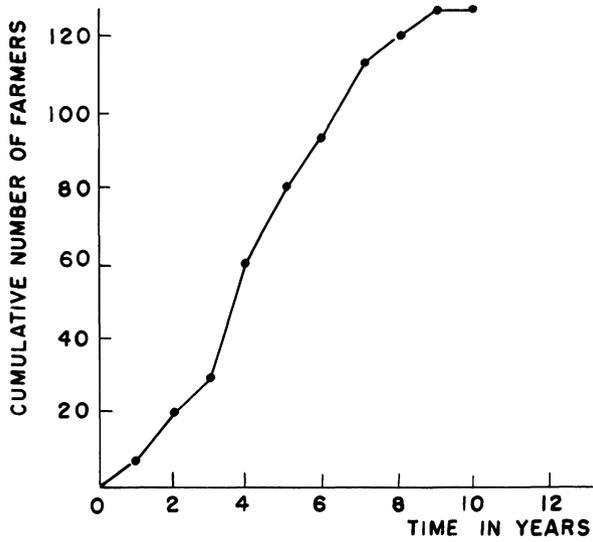


Figure 1. Cumulative number of farmers who adopted 2,4-D weed spray in Iowa

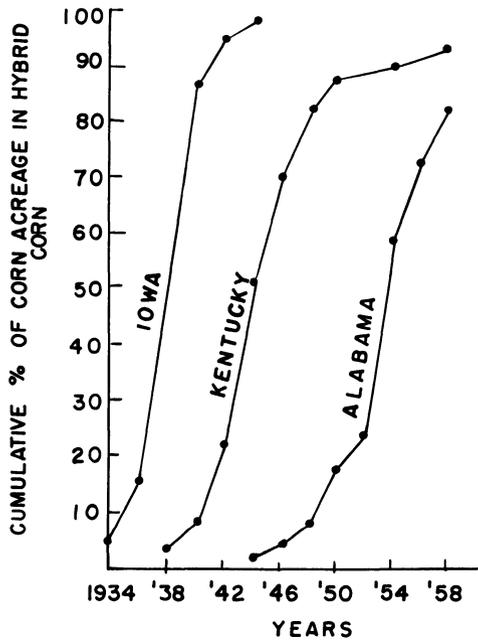


Figure 2. Cumulative percentage of corn acreage in hybrid corn in three American states

as radio, television, newspapers and farmers' magazines play a large role in the early stages of the adoption process. Therefore, we must add a term to the differential equation (1) to take this into account. To compute this term, we assume that the number of farmers Δp who learn of the innovation through the mass communication media in a short period of time Δt is proportional to the number of farmers who do not yet know; i.e.,

$$\Delta p = c'(N - p)\Delta t$$

for some positive constant c' . Letting $\Delta t \rightarrow 0$, we see that $c'(N - p)$ farmers, per unit time, learn of the innovation through the mass communication media. Thus, if $p(0) = 0$, then $p(t)$ satisfies the initial-value problem

$$\frac{dp}{dt} = cp(N - p) + c'(N - p), \quad p(0) = 0. \quad (4)$$

The solution of (4) is

$$p(t) = \frac{Nc' [e^{(c' + cN)t} - 1]}{cN + c'e^{(c' + cN)t}}, \quad (5)$$

and in Exercises 2 and 3 we indicate how to determine the shape of the curve (5).

The corrected curve (5) now gives remarkably good agreement with Figures 1 and 2, for suitable choices of c and c' . However, (see Exercise 3c) it still doesn't explain why the adoption of hybrid corn in Alabama only began to decelerate after sixty per cent of the farmers had adopted the innovation. This indicates, of course, that other factors, such as the time interval that elapses between when a farmer first learns of an innovation and when he actually adopts it, may play an important role in the adoption process, and must be taken into account in any model.

We would now like to show that the differential equation

$$dp/dt = cp(N - p)$$

also governs the rate at which firms in such diverse industries as bituminous coal, iron and steel, brewing, and railroads adopted several major innovations in the first part of this century. This is rather surprising, since we would expect that the number of firms adopting an innovation in one of these industries certainly depends on the profitability of the innovation and the investment required to implement it, and we haven't mentioned these factors in deriving Equation (1). However, as we shall see shortly, these two factors are incorporated in the constant c .

Let n be the total number of firms in a particular industry who have adopted an innovation at time t . It is clear that the number of firms Δp who adopt the innovation in a short time interval Δt is proportional to the number of firms $n - p$ who have not yet adopted; i.e., $\Delta p = \lambda(n - p)\Delta t$. Letting $\Delta t \rightarrow 0$, we see that

$$\frac{dp}{dt} = \lambda(n - p).$$

The proportionality factor λ depends on the profitability π of installing this innovation relative to that of alternative investments, the investment s required to install this innovation as a percentage of the total assets of the firm, and the percentage of firms who have already adopted. Thus,

$$\lambda = f(\pi, s, p/n).$$

Expanding f in a Taylor series, and dropping terms of degree more than two, gives

$$\begin{aligned} \lambda = & a_1 + a_2\pi + a_3s + a_4\frac{p}{n} + a_5\pi^2 + a_6s^2 + a_7\pi s \\ & + a_8\pi\left(\frac{p}{n}\right) + a_9s\left(\frac{p}{n}\right) + a_{10}\left(\frac{p}{n}\right)^2. \end{aligned}$$

In the late 1950's, Edwin Mansfield of Carnegie Mellon University investigated the spread of twelve innovations in four major industries. From his exhaustive studies, Mansfield concluded that $a_{10} = 0$ and

$$a_1 + a_2\pi + a_3s + a_5\pi^2 + a_6s^2 + a_7\pi s = 0.$$

Thus, setting

$$k = a_4 + a_8\pi + a_9s, \quad (6)$$

we see that

$$\frac{dp}{dt} = k\frac{p}{n}(n-p).$$

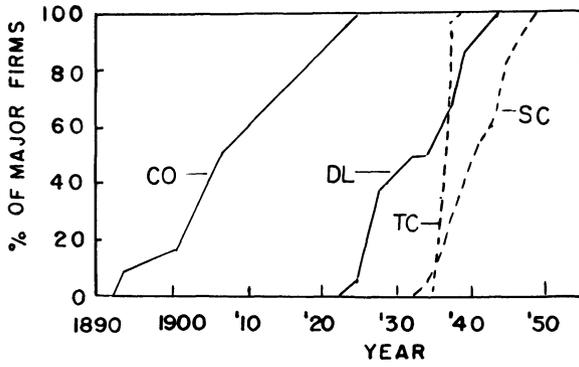
(This is the equation obtained previously for the spread of innovations among farmers, if we set $k/n = c$.) We assume that the innovation is first adopted by one firm in the year t_0 . Then, $p(t)$ satisfies the initial-value problem

$$\frac{dp}{dt} = \frac{k}{n}p(n-p), \quad p(t_0) = 1 \quad (7)$$

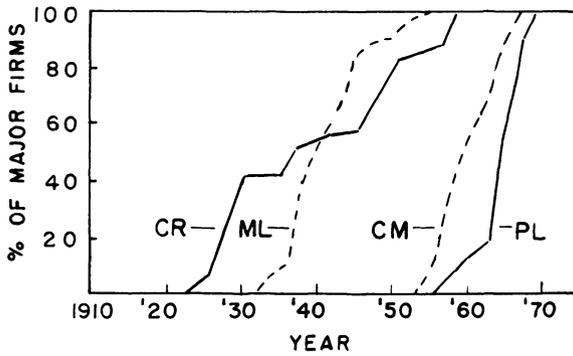
and this implies that

$$p(t) = \frac{n}{1 + (n-1)e^{-k(t-t_0)}}.$$

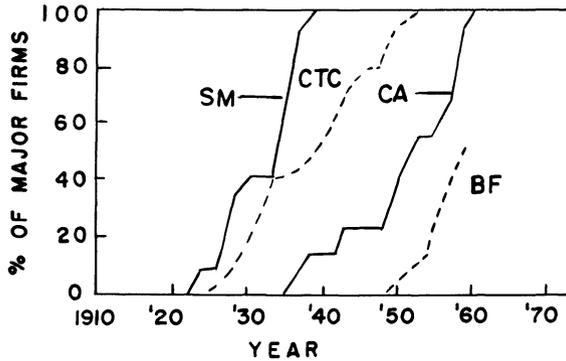
Mansfield studied how rapidly the use of twelve innovations spread from enterprise to enterprise in four major industries—bituminous coal, iron and steel, brewing, and railroads. The innovations are the shuttle car, trackless mobile loader, and continuous mining machine (in bituminous coal); the by-product coke oven, continuous wide strip mill, and continuous annealing line for tin plate (in iron and steel); the pallet-loading machine, tin container, and high speed bottle filler (in brewing); and the diesel locomotive, centralized traffic control, and car retarders (in railroads). His results are described graphically in Figure 3. For all but the by-product coke oven and tin container, the percentages given are for every two years from the year of initial introduction. The length of the interval for the by-product coke oven is about six years, and for the tin container, it is six months. Notice how all these curves have the general appearance of a logistic curve.



(a)



(b)



(c)

Figure 3. Growth in the percentage of major firms that introduced twelve innovations; bituminous coal, iron and steel, brewing, and railroad industries, 1890–1958; (a) By-product coke oven (CO), diesel locomotive (DL), tin container (TC), and shuttle car (SC); (b) Car retarder (CR), trackless mobile loader (ML), continuous mining machine (CM), and pallet-loading machine (PL); (c) Continuous wide-strip mill (SM), centralized traffic control (CTC), continuous annealing (CA), and highspeed bottle filler (BF).

Table 1.

Innovation	n	t_0	a_4	a_8	a_9	π	s
Diesel locomotive	25	1925	-0.59	0.530	-0.027	1.59	0.015
Centralized traffic control	24	1926	-0.59	0.530	-0.027	1.48	0.024
Car retarders	25	1924	-0.59	0.530	-0.027	1.25	0.785
Continuous wide strip mill	12	1924	-0.52	0.530	-0.027	1.87	4.908
By-product coke oven	12	1894	-0.52	0.530	-0.027	1.47	2.083
Continuous annealing	9	1936	-0.52	0.530	-0.027	1.25	0.554
Shuttle car	15	1937	-0.57	0.530	-0.027	1.74	0.013
Trackless mobile loader	15	1934	-0.57	0.530	-0.027	1.65	0.019
Continuous mining machine	17	1947	-0.57	0.530	-0.027	2.00	0.301
Tin container	22	1935	-0.29	0.530	-0.027	5.07	0.267
High speed bottle filler	16	1951	-0.29	0.530	-0.027	1.20	0.575
Pallet-loading machine	19	1948	-0.29	0.530	-0.027	1.67	0.115

For a more detailed comparison of the predictions of our model (7) with these observed results, we must evaluate the constants n , k , and t_0 for each of the twelve innovations. Table 1 gives the value of n , t_0 , a_4 , a_5 , a_9 , π , and s for each of the twelve innovations; the constant k can then be computed from Equation (6). As the answers to Exercises 5 and 6 will indicate, our model (7) predicts with reasonable accuracy the rate of adoption of these twelve innovations.

Reference

Mansfield, E., "Technical change and the rate of imitation," *Econometrica*, Vol. 29, No. 4, Oct. 1961.

EXERCISES

1. Solve the initial-value problem (2).
2. Let $c=0$ in (5). Show that $p(t)$ increases monotonically from 0 to N , and has no points of inflection.
3. Here is a heuristic argument to determine the behavior of the curve (5). If $c'=0$, then we have a logistic curve, and if $c=0$, then we have the behavior described in Exercise 2. Thus, if c is large relative to c' , then we have a logistic curve, and if c is small relative to c' then we have the behavior illustrated in Exercise 2.

1 First-order differential equations

- (a) Let $p(t)$ satisfy (4). Show that

$$\frac{d^2p}{dt^2} = (N - p)(cp + c')(cN - 2cp - c').$$

- (b) Show that $p(t)$ has a point of inflection, at which dp/dt achieves a maximum, if, and only if, $c'/c < N$.
- (c) Assume that $p(t)$ has a point of inflection at $t = t^*$. Show that $p(t^*) \leq N/2$.
4. Solve the initial-value problem (7).
5. It seems reasonable to take the time span between the date when 20% of the firms had introduced the innovation and the date when 80% of the firms had introduced the innovation, as the rate of imitation.
- (a) Show from our model that this time span is $4(\ln 2)/k$.
- (b) For each of the twelve innovations, compute this time span from the data in Table 1, and compare with the observed value in Figure 3.
6. (a) Show from our model that $(1/k)\ln(n - 1)$ years elapse before 50% of the firms introduce an innovation.
- (b) Compute this time span for each of the 12 innovations and compare with the observed values in Figure 3.

1.7 An atomic waste disposal problem

For several years the Atomic Energy Commission (now known as the Nuclear Regulatory Commission) had disposed of concentrated radioactive waste material by placing it in tightly sealed drums which were then dumped at sea in fifty fathoms (300 feet) of water. When concerned ecologists and scientists questioned this practice, they were assured by the A.E.C. that the drums would never develop leaks. Exhaustive tests on the drums proved the A.E.C. right. However, several engineers then raised the question of whether the drums could crack from the impact of hitting the ocean floor. "Never," said the A.E.C. "We'll see about that," said the engineers. After performing numerous experiments, the engineers found that the drums could crack on impact if their velocity exceeded forty feet per second. The problem before us, therefore, is to compute the velocity of the drums upon impact with the ocean floor. To this end, we digress briefly to study elementary Newtonian mechanics.

Newtonian mechanics is the study of Newton's famous laws of motion and their consequences. Newton's first law of motion states that an object will remain at rest, or move with constant velocity, if no force is acting on it. A force should be thought of as a push or pull. This push or pull can be exerted directly by something in contact with the object, or it can be exerted indirectly, as the earth's pull of gravity is.

Newton's second law of motion is concerned with describing the motion of an object which is acted upon by several forces. Let $y(t)$ denote the position of the center of gravity of the object. (We assume that the object moves in only one direction.) Those forces acting on the object, which tend

to increase y , are considered positive, while those forces tending to decrease y are considered negative. The resultant force F acting on an object is defined to be the sum of all positive forces minus the sum of all negative forces. Newton's second law of motion states that the acceleration d^2y/dt^2 of an object is proportional to the resultant force F acting on it; i.e.,

$$\frac{d^2y}{dt^2} = \frac{1}{m} F. \quad (1)$$

The constant m is the mass of the object. It is related to the weight W of the object by the relation $W = mg$, where g is the acceleration of gravity. Unless otherwise stated, we assume that the weight of an object and the acceleration of gravity are constant. We will also adopt the English system of units, so that t is measured in seconds, y is measured in feet, and F is measured in pounds. The units of m are then slugs, and the gravitational acceleration g equals 32.2 ft/s^2 .

Remark. We would prefer to use the mks system of units, where y is measured in meters and F is measured in newtons. The units of m are then kilograms, and the gravitational acceleration equals 9.8 m/s^2 . In the third edition of this text, we have changed from the English system of units to the mks system in Section 2.6. However, changing to the mks system in this section would have caused undue confusion to the users of the first and second editions. This is because of the truncation error involved in converting from feet to meters and pounds to newtons.

We return now to our atomic waste disposal problem. As a drum descends through the water, it is acted upon by three forces W , B , and D . The force W is the weight of the drum pulling it down, and in magnitude, $W = 527.436 \text{ lb}$. The force B is the buoyancy force of the water acting on the drum. This force pushes the drum up, and its magnitude is the weight of the water displaced by the drum. Now, the Atomic Energy Commission used 55 gallon drums, whose volume is 7.35 ft^3 . The weight of one cubic foot of salt water is 63.99 lb . Hence $B = (63.99)(7.35) = 470.327 \text{ lb}$.

The force D is the drag force of the water acting on the drum; it resists the motion of the drum through the water. Experiments have shown that any medium such as water, oil, and air resists the motion of an object through it. This resisting force acts in the direction opposite the motion, and is usually directly proportional to the velocity V of the object. Thus, $D = cV$, for some positive constant c . Notice that the drag force increases as V increases, and decreases as V decreases. To calculate D , the engineers conducted numerous towing experiments. They concluded that the orientation of the drum had little effect on the drag force, and that

$$D = 0.08 V \frac{(\text{lb})(\text{s})}{\text{ft}}.$$

1 First-order differential equations

Now, set $y=0$ at sea level, and let the direction of increasing y be downwards. Then, W is a positive force, and B and D are negative forces. Consequently, from (1),

$$\frac{d^2y}{dt^2} = \frac{1}{m}(W - B - cV) = \frac{g}{W}(W - B - cV).$$

We can rewrite this equation as a first-order linear differential equation for $V = dy/dt$; i.e.,

$$\frac{dV}{dt} + \frac{cg}{W}V = \frac{g}{W}(W - B). \quad (2)$$

Initially, when the drum is released in the ocean, its velocity is zero. Thus, $V(t)$, the velocity of the drum, satisfies the initial-value problem

$$\frac{dV}{dt} + \frac{cg}{W}V = \frac{g}{W}(W - B), \quad V(0) = 0, \quad (3)$$

and this implies that

$$V(t) = \frac{W - B}{c} [1 - e^{(-cg/W)t}]. \quad (4)$$

Equation (4) expresses the velocity of the drum as a function of time. In order to determine the impact velocity of the drum, we must compute the time t at which the drum hits the ocean floor. Unfortunately, though, it is impossible to find t as an explicit function of y (see Exercise 2). Therefore, we cannot use Equation (4) to find the velocity of the drum when it hits the ocean floor. However, the A.E.C. can use this equation to try and prove that the drums do not crack on impact. To wit, observe from (4) that $V(t)$ is a monotonic increasing function of time which approaches the limiting value

$$V_T = \frac{W - B}{c}$$

as t approaches infinity. The quantity V_T is called the terminal velocity of the drum. Clearly, $V(t) \leq V_T$, so that the velocity of the drum when it hits the ocean floor is certainly less than $(W - B)/c$. Now, if this terminal velocity is less than 40 ft/s, then the drums could not possibly break on impact. However,

$$\frac{W - B}{c} = \frac{527.436 - 470.327}{0.08} = 713.86 \text{ ft/s,}$$

and this is way too large.

It should be clear now that the only way we can resolve the dispute between the A.E.C. and the engineers is to find $v(y)$, the velocity of the drum as a function of position. The function $v(y)$ is very different from the function $V(t)$, which is the velocity of the drum as a function of time. However, these two functions are related through the equation

$$V(t) = v(y(t))$$

if we express y as a function of t . By the chain rule of differentiation, $dV/dt = (dv/dy)(dy/dt)$. Hence

$$\frac{W}{g} \frac{dv}{dy} \frac{dy}{dt} = W - B - cV.$$

But $dy/dt = V(t) = v(y(t))$. Thus, suppressing the dependence of y on t , we see that $v(y)$ satisfies the first-order differential equation

$$\frac{W}{g} v \frac{dv}{dy} = W - B - cv, \quad \text{or} \quad \frac{v}{W - B - cv} \frac{dv}{dy} = \frac{g}{W}.$$

Moreover,

$$v(0) = v(y(0)) = V(0) = 0.$$

Hence,

$$\int_0^v \frac{r \, dr}{W - B - cr} = \int_0^y \frac{g}{W} \, ds = \frac{gy}{W}.$$

Now,

$$\begin{aligned} \int_0^v \frac{r \, dr}{W - B - cr} &= \int_0^v \frac{r - (W - B)/c}{W - B - cr} \, dr + \frac{W - B}{c} \int_0^v \frac{dr}{W - B - cr} \\ &= -\frac{1}{c} \int_0^v dr + \frac{W - B}{c} \int_0^v \frac{dr}{W - B - cr} \\ &= -\frac{v}{c} - \frac{(W - B)}{c^2} \ln \frac{|W - B - cv|}{W - B}. \end{aligned}$$

We know already that $v < (W - B)/c$. Consequently, $W - B - cv$ is always positive, and

$$\frac{gy}{W} = -\frac{v}{c} - \frac{(W - B)}{c^2} \ln \frac{W - B - cv}{W - B}. \quad (5)$$

At this point, we are ready to scream in despair since we cannot find v as an explicit function of y from (5). This is not an insurmountable difficulty, though. As we show in Section 1.11, it is quite simple, with the aid of a digital computer, to find $v(300)$ from (5). We need only supply the computer with a good approximation of $v(300)$ and this is obtained in the following manner. The velocity $v(y)$ of the drum satisfies the initial-value problem

$$\frac{W}{g} v \frac{dv}{dy} = W - B - cv, \quad v(0) = 0. \quad (6)$$

Let us, for the moment, set $c = 0$ in (6) to obtain the new initial-value problem

$$\frac{W}{g} u \frac{du}{dy} = W - B, \quad u(0) = 0. \quad (6')$$

(We have replaced v by u to avoid confusion later.) We can integrate (6') immediately to obtain that

$$\frac{W}{g} \frac{u^2}{2} = (W - B)y, \quad \text{or} \quad u(y) = \left[\frac{2g}{W} (W - B)y \right]^{1/2}.$$

In particular,

$$u(300) = \left[\frac{2g}{W} (W - B) 300 \right]^{1/2} = \left[\frac{2(32.2)(57.109)(300)}{527.436} \right]^{1/2} \\ \cong \sqrt{2092} \cong 45.7 \text{ ft/s.}$$

We claim, now, that $u(300)$ is a very good approximation of $v(300)$. The proof of this is as follows. First, observe that the velocity of the drum is always greater if there is no drag force opposing the motion. Hence,

$$v(300) < u(300).$$

Second, the velocity v increases as y increases, so that $v(y) \leq v(300)$ for $y \leq 300$. Consequently, the drag force D of the water acting on the drum is always less than $0.08 \times u(300) \cong 3.7$ lb. Now, the resultant force $W - B$ pulling the drum down is approximately 57.1 lb, which is very large compared to D . It stands to reason, therefore, that $u(y)$ should be a very good approximation of $v(y)$. And indeed, this is the case, since we find numerically (see Section 1.11) that $v(300) = 45.1$ ft/s. Thus, the drums can break upon impact, and the engineers were right.

Epilog. The rules of the Atomic Energy Commission now expressly forbid the dumping of low level atomic waste at sea. This author is uncertain though, as to whether Western Europe has also forbidden this practice.

Remark. The methods introduced in this section can also be used to find the velocity of any object which is moving through a medium that resists the motion. We just disregard the buoyancy force if the medium is not water. For example, let $V(t)$ denote the velocity of a parachutist falling to earth under the influence of gravity. Then,

$$\frac{W}{g} \frac{dV}{dt} = W - D$$

where W is the weight of the man and the parachute, and D is the drag force exerted by the atmosphere on the falling parachutist. The drag force on a bluff object in air, or in any fluid of small viscosity is usually very nearly proportional to V^2 . Proportionality to V is the exceptional case, and occurs only at very low speeds. The criterion as to whether the square or the linear law applies is the "Reynolds number"

$$R = \rho V L / \mu.$$

L is a representative length dimension of the object, and ρ and μ are the density and viscosity of the fluid. If $R < 10$, then $D \sim V$, and if $R > 10^3$, $D \sim V^2$. For $10 < R < 10^3$, neither law is accurate.

EXERCISES

1. Solve the initial-value problem (3).
2. Solve for $y = y(t)$ from (4), and then show that the equation $y = y(t)$ cannot be solved explicitly for $t = t(y)$.
3. Show that the drums of atomic waste will not crack upon impact if they are dropped into L feet of water with $(2g(W - B)L/W)^{1/2} < 40$.
4. Fat Richie, an enormous underworld hoodlum weighing 400 lb, was pushed out of a penthouse window 2800 feet above the ground in New York City. Neglecting air resistance find (a) the velocity with which Fat Richie hit the ground; (b) the time elapsed before Fat Richie hit the ground.
5. An object weighing 300 lb is dropped into a river 150 feet deep. The volume of the object is 2 ft^3 , and the drag force exerted by the water on it is 0.05 times its velocity. The drag force may be considered negligible if it does not exceed 5% of the resultant force pulling the drum down. Prove that the drag force is negligible in this case. (Here $B = 2(62.4) = 124.8$.)
6. A 400 lb sphere of volume $4\pi/3$ and a 300 lb cylinder of volume π are simultaneously released from rest into a river. The drag force exerted by the water on the falling sphere and cylinder is λV_s and λV_c , respectively, where V_s and V_c are the velocities of the sphere and cylinder, and λ is a positive constant. Determine which object reaches the bottom of the river first.
7. A parachutist falls from rest toward earth. The combined weight of man and parachute is 161 lb. Before the parachute opens, the air resistance equals $V/2$. The parachute opens 5 seconds after the fall begins; and the air resistance is then $V^2/2$. Find the velocity $V(t)$ of the parachutist after the parachute opens.
8. A man wearing a parachute jumps from a great height. The combined weight of man and parachute is 161 lb. Let $V(t)$ denote his speed at time t seconds after the fall begins. During the first 10 seconds, the air resistance is $V/2$. Thereafter, while the parachute is open, the air resistance is $10V$. Find an explicit formula for $V(t)$ at any time t greater than 10 seconds.
9. An object of mass m is projected vertically downward with initial velocity V_0 in a medium offering resistance proportional to the square root of the magnitude of the velocity.
 - (a) Find a relation between the velocity V and the time t if the drag force equals $c\sqrt{V}$.
 - (b) Find the terminal velocity of the object. *Hint*: You can find the terminal velocity even though you cannot solve for $V(t)$.
10. A body of mass m falls from rest in a medium offering resistance proportional to the square of the velocity; that is, $D = cV^2$. Find $V(t)$ and compute the terminal velocity V_T .
11. A body of mass m is projected upward from the earth's surface with an initial velocity V_0 . Take the y -axis to be positive upward, with the origin on the surface of the earth. Assuming there is no air resistance, but taking into

account the variation of the earth's gravitational field with altitude, we obtain that

$$m \frac{dV}{dt} = - \frac{mgR^2}{(y+R)^2}$$

where R is the radius of the earth.

- (a) Let $V(t) = v(y(t))$. Find a differential equation satisfied by $v(y)$.
 (b) Find the smallest initial velocity V_0 for which the body will not return to earth. This is the so-called escape velocity. *Hint*: The escape velocity is found by requiring that $v(y)$ remain strictly positive.
12. It is not really necessary to find $v(y)$ explicitly in order to prove that $v(300)$ exceeds 40 ft/s. Here is an alternate proof. Observe first that $v(y)$ increases as y increases. This implies that y is a monotonic increasing function of v . Therefore, if y is less than 300 ft when v is 40 ft/s, then v must be greater than 40 ft/s when y is 300 ft. Substitute $v = 40$ ft/s in Equation (5), and show that y is less than 300 ft. Conclude, therefore, that the drums can break upon impact.

1.8 The dynamics of tumor growth, mixing problems and orthogonal trajectories

In this section we present three very simple but extremely useful applications of first-order equations. The first application concerns the growth of solid tumors; the second application is concerned with "mixing problems" or "compartment analysis"; and the third application shows how to find a family of curves which is orthogonal to a given family of curves.

(a) *The dynamics of tumor growth*

It has been observed experimentally, that "free living" dividing cells, such as bacteria cells, grow at a rate proportional to the volume of dividing cells at that moment. Let $V(t)$ denote the volume of dividing cells at time t . Then,

$$\frac{dV}{dt} = \lambda V \tag{1}$$

for some positive constant λ . The solution of (1) is

$$V(t) = V_0 e^{\lambda(t-t_0)} \tag{2}$$

where V_0 is the volume of dividing cells at the initial time t_0 . Thus, free living dividing cells grow *exponentially* with time. One important consequence of (2) is that the volume of cells keeps doubling (see Exercise 1) every time interval of length $\ln 2/\lambda$.

On the other hand, solid tumors do not grow exponentially with time. As the tumor becomes larger, the doubling time of the total tumor volume continuously increases. Various researchers have shown that the data for many solid tumors is fitted remarkably well, over almost a 1000 fold in-

crease in tumor volume, by the equation

$$V(t) = V_0 \exp\left(\frac{\lambda}{\alpha}(1 - \exp(-\alpha t))\right) \quad (3)$$

where $\exp(x) \equiv e^x$, and λ and α are positive constants.

Equation (3) is usually referred to as a Gompertzian relation. It says that the tumor grows more and more slowly with the passage of time, and that it ultimately approaches the limiting volume $V_0 e^{\lambda/\alpha}$. Medical researchers have long been concerned with explaining this deviation from simple exponential growth. A great deal of insight into this problem can be gained by finding a differential equation satisfied by $V(t)$. Differentiating (3) gives

$$\begin{aligned} \frac{dV}{dt} &= V_0 \lambda \exp(-\alpha t) \exp\left(\frac{\lambda}{\alpha}(1 - \exp(-\alpha t))\right) \\ &= \lambda e^{-\alpha t} V. \end{aligned} \quad (4)$$

Two conflicting theories have been advanced for the dynamics of tumor growth. They correspond to the two arrangements

$$\frac{dV}{dt} = (\lambda e^{-\alpha t}) V \quad (4a)$$

$$\frac{dV}{dt} = \lambda(e^{-\alpha t} V) \quad (4b)$$

of the differential equation (4). According to the first theory, the retarding effect of tumor growth is due to an increase in the mean generation time of the cells, without a change in the proportion of reproducing cells. As time goes on, the reproducing cells mature, or age, and thus divide more slowly. This theory corresponds to the bracketing (a).

The bracketing (b) suggests that the mean generation time of the dividing cells remains constant, and the retardation of growth is due to a loss in reproductive cells in the tumor. One possible explanation for this is that a *necrotic region* develops in the center of the tumor. This necrosis appears at a critical size for a particular type of tumor, and thereafter the necrotic "core" increases rapidly as the total tumor mass increases. According to this theory, a necrotic core develops because in many tumors the supply of blood, and thus of oxygen and nutrients, is almost completely confined to the surface of the tumor and a short distance beneath it. As the tumor grows, the supply of oxygen to the central core by diffusion becomes more and more difficult resulting in the formation of a necrotic core.

(b) *Mixing problems*

Many important problems in biology and engineering can be put into the following framework. A solution containing a fixed concentration of substance x flows into a tank, or compartment, containing the substance x and possibly other substances, at a specified rate. The mixture is stirred

1 First-order differential equations

together very rapidly, and then leaves the tank, again at a specified rate. Find the concentration of substance x in the tank at any time t .

Problems of this type fall under the general heading of “mixing problems,” or compartment analysis. The following example illustrates how to solve these problems.

Example 1. A tank contains S_0 lb of salt dissolved in 200 gallons of water. Starting at time $t=0$, water containing $\frac{1}{2}$ lb of salt per gallon enters the tank at the rate of 4 gal/min, and the well stirred solution leaves the tank at the same rate. Find the concentration of salt in the tank at any time $t > 0$.

Solution. Let $S(t)$ denote the amount of salt in the tank at time t . Then, $S'(t)$, which is the rate of change of salt in the tank at time t , must equal the rate at which salt enters the tank minus the rate at which it leaves the tank. Obviously, the rate at which salt enters the tank is

$$\frac{1}{2} \text{ lb/gal times } 4 \text{ gal/min} = 2 \text{ lb/min.}$$

After a moment's reflection, it is also obvious that the rate at which salt leaves the tank is

$$4 \text{ gal/min times } \frac{S(t)}{200}.$$

Thus

$$S'(t) = 2 - \frac{S(t)}{50}, \quad S(0) = S_0,$$

and this implies that

$$S(t) = S_0 e^{-0.02t} + 100(1 - e^{-0.02t}). \quad (5)$$

Hence, the concentration $c(t)$ of salt in the tank is given by

$$c(t) = \frac{S(t)}{200} = \frac{S_0}{200} e^{-0.02t} + \frac{1}{2}(1 - e^{-0.02t}). \quad (6)$$

Remark. The first term on the right-hand side of (5) represents the portion of the original amount of salt remaining in the tank at time t . This term becomes smaller and smaller with the passage of time as the original solution is drained from the tank. The second term on the right-hand side of (5) represents the amount of salt in the tank at time t due to the action of the flow process. Clearly, the amount of salt in the tank must ultimately approach the limiting value of 100 lb, and this is easily verified by letting t approach ∞ in (5).

(c) *Orthogonal trajectories*

In many physical applications, it is often necessary to find the orthogonal trajectories of a given family of curves. (A curve which intersects each

member of a family of curves at right angles is called an orthogonal trajectory of the given family.) For example, a charged particle moving under the influence of a magnetic field always travels on a curve which is perpendicular to each of the magnetic field lines. The problem of computing orthogonal trajectories of a family of curves can be solved in the following manner. Let the given family of curves be described by the relation

$$F(x, y, c) = 0. \quad (7)$$

Differentiating this equation yields

$$F_x + F_y y' = 0, \quad \text{or} \quad y' = -\frac{F_x}{F_y}. \quad (8)$$

Next, we solve for $c = c(x, y)$ from (7) and replace every c in (8) by this value $c(x, y)$. Finally, since the slopes of curves which intersect orthogonally are negative reciprocals of each other, we see that the orthogonal trajectories of (7) are the solution curves of the equation

$$y' = \frac{F_y}{F_x}. \quad (9)$$

Example 2. Find the orthogonal trajectories of the family of parabolas

$$x = cy^2.$$

Solution. Differentiating the equation $x = cy^2$ gives $1 = 2cyy'$. Since $c = x/y^2$, we see that $y' = y/2x$. Thus, the orthogonal trajectories of the family

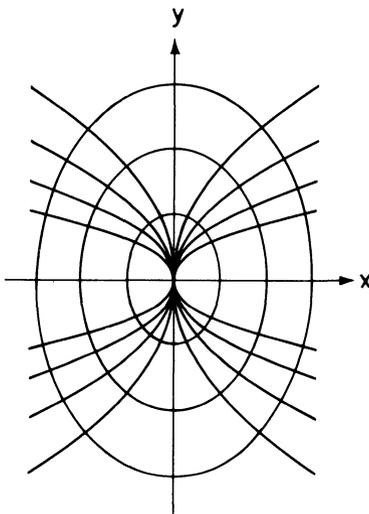


Figure 1. The parabolas $x = cy^2$ and their orthogonal trajectories

1 First-order differential equations

of parabolas $x = cy^2$ are the solution curves of the equation

$$y' = -\frac{2x}{y}. \quad (10)$$

This equation is separable, and its solution is

$$y^2 + 2x^2 = k^2. \quad (11)$$

Thus, the family of ellipses (11) (see Figure 1) are the orthogonal trajectories of the family of parabolas $x = cy^2$.

Reference

Burton, Alan C., Rate of growth of solid tumors as a problem of diffusion, *Growth*, 1966, vol. 30, pp. 157–176.

EXERCISES

1. A given substance satisfies the exponential growth law (1). Show that the graph of $\ln V$ versus t is a straight line.
2. A substance x multiplies exponentially, and a given quantity of the substance doubles every 20 years. If we have 3 lb of substance x at the present time, how many lb will we have 7 years from now?
3. A substance x decays exponentially, and only half of the given quantity of x remains after 2 years. How long does it take for 5 lb of x to decay to 1 lb?
4. The equation $p' = ap^\alpha$, $\alpha > 1$, is proposed as a model of the population growth of a certain species. Show that $p(t) \rightarrow \infty$ in finite time. Conclude, therefore, that this model is not accurate over a reasonable length of time.
5. A cancerous tumor satisfies the Gompertzian relation (3). Originally, when it contained 10^4 cells, the tumor was increasing at the rate of 20% per unit time. The numerical value of the retarding constant α is 0.02. What is the limiting number of cells in this tumor?
6. A tracer dose of radioactive iodine ^{131}I is injected into the blood stream at time $t = 0$. Assume that the original amount Q_0 of iodine is distributed evenly in the entire blood stream before any loss occurs. Let $Q(t)$ denote the amount of iodine in the blood at time $t > 0$. Part of the iodine leaves the blood and enters the urine at the rate $k_1 Q$. Another part of the iodine enters the thyroid gland at the rate $k_2 Q$. Find $Q(t)$.
7. Industrial waste is pumped into a tank containing 1000 gallons of water at the rate of 1 gal/min, and the well-stirred mixture leaves the tank at the same rate. (a) Find the concentration of waste in the tank at time t . (b) How long does it take for the concentration to reach 20%?
8. A tank contains 300 gallons of water and 100 gallons of pollutants. Fresh water is pumped into the tank at the rate of 2 gal/min, and the well-stirred mixture leaves at the same rate. How long does it take for the concentration of pollutants in the tank to decrease to 1/10 of its original value?

9. Consider a tank containing, at time $t=0$, Q_0 lb of salt dissolved in 150 gallons of water. Assume that water containing $\frac{1}{2}$ lb of salt per gallon is entering the tank at a rate of 3 gal/min, and that the well-stirred solution is leaving the tank at the same rate. Find an expression for the concentration of salt in the tank at time t .
10. A room containing 1000 cubic feet of air is originally free of carbon monoxide. Beginning at time $t=0$, cigarette smoke containing 4 percent carbon monoxide is blown into the room at the rate of $0.1 \text{ ft}^3/\text{min}$, and the well-circulated mixture leaves the room at the same rate. Find the time when the concentration of carbon monoxide in the room reaches 0.012 percent. (Extended exposure to this concentration of carbon monoxide is dangerous.)
11. A 500 gallon tank originally contains 100 gallons of fresh water. Beginning at time $t=0$, water containing 50 percent pollutants flows into the tank at the rate of 2 gal/min, and the well-stirred mixture leaves at the rate of 1 gal/min. Find the concentration of pollutants in the tank at the moment it overflows.

In Exercises 12–17, find the orthogonal trajectories of the given family of curves.

12. $y = cx^2$

13. $y^2 - x^2 = c$

14. $y = c \sin x$

15. $x^2 + y^2 = cx$ (see Exercise 13 of Section 1.4)

16. $y = ce^x$

17. $y = e^{cx}$

18. The presence of toxins in a certain medium destroys a strain of bacteria at a rate jointly proportional to the number of bacteria present and to the amount of toxin. Call the constant of proportionality a . If there were no toxins present, the bacteria would grow at a rate proportional to the amount present. Call this constant of proportionality b . Assume that the amount T of toxin is increasing at a constant rate c , that is, $dT/dt = c$, and that the production of toxins begins at time $t=0$. Let $y(t)$ denote the number of living bacteria present at time t .
 - (a) Find a first-order differential equation satisfied by $y(t)$.
 - (b) Solve this differential equation to obtain $y(t)$. What happens to $y(t)$ as t approaches ∞ ?
19. Many savings banks now advertise continuous compounding of interest. This means that the amount of money $P(t)$ on deposit at time t , satisfies the differential equation $dP(t)/dt = rP(t)$ where r is the annual interest rate and t is measured in years. Let P_0 denote the original principal.
 - (a) Show that $P(1) = P_0 e^r$.
 - (b) Let $r = 0.0575, 0.065, 0.0675, \text{ and } 0.075$. Show that $e^r = 1.05919, 1.06716, 1.06983, \text{ and } 1.07788$, respectively. Thus, the effective annual yield on interest rates of $5\frac{3}{4}\%$, $6\frac{1}{2}\%$, $6\frac{3}{4}\%$, and $7\frac{1}{2}\%$ should be 5.919, 6.716, 6.983, and 7.788%, respectively. Most banks, however, advertise effective annual yields of 6, 6.81, 7.08, and 7.9%, respectively. The reason for this discrepancy is that banks calculate a daily rate of interest based on 360 days, and they pay interest for each day money is on deposit. For a year, one gets five extra

days. Thus, we must multiply the annual yields of 5.919, 6.716, 6.983, and 7.788% by $365/360$, and then we obtain the advertised values.

- (c) It is interesting to note that the Old Colony Cooperative Bank in Rhode Island advertises an effective annual yield of 6.72% on an annual interest rate of $6\frac{1}{2}\%$ (the lower value), and an effective annual yield of 7.9% on an annual interest rate of $7\frac{1}{2}\%$. Thus they are inconsistent.

1.9 Exact equations, and why we cannot solve very many differential equations

When we began our study of differential equations, the only equation we could solve was $dy/dt = g(t)$. We then enlarged our inventory to include all linear and separable equations. More generally, we can solve all differential equations which are, or can be put, in the form

$$\frac{d}{dt}\phi(t,y) = 0 \quad (1)$$

for some function $\phi(t,y)$. To wit, we can integrate both sides of (1) to obtain that

$$\phi(t,y) = \text{constant} \quad (2)$$

and then solve for y as a function of t from (2).

Example 1. The equation $1 + \cos(t+y) + \cos(t+y)(dy/dt) = 0$ can be written in the form $(d/dt)[t + \sin(t+y)] = 0$. Hence,

$$\phi(t,y) = t + \sin(t+y) = c, \quad \text{and} \quad y = -t + \arcsin(c-t).$$

Example 2. The equation $\cos(t+y) + [1 + \cos(t+y)]dy/dt = 0$ can be written in the form $(d/dt)[y + \sin(t+y)] = 0$. Hence,

$$\phi(t,y) = y + \sin(t+y) = c.$$

We must leave the solution in this form though, since we cannot solve for y explicitly as a function of time.

Equation (1) is clearly the most general first-order differential equation that we can solve. Thus, it is important for us to be able to recognize when a differential equation can be put in this form. This is not as simple as one might expect. For example, it is certainly not obvious that the differential equation

$$2t + y - \sin t + (3y^2 + \cos y + t) \frac{dy}{dt} = 0$$

can be written in the form $(d/dt)(y^3 + t^2 + ty + \sin y + \cos t) = 0$. To find all those differential equations which can be written in the form (1), observe,

from the chain rule of partial differentiation, that

$$\frac{d}{dt}\phi(t,y(t)) = \frac{\partial\phi}{\partial t} + \frac{\partial\phi}{\partial y} \frac{dy}{dt}.$$

Hence, the differential equation $M(t,y) + N(t,y)(dy/dt) = 0$ can be written in the form $(d/dt)\phi(t,y) = 0$ if and only if there exists a function $\phi(t,y)$ such that $M(t,y) = \partial\phi/\partial t$ and $N(t,y) = \partial\phi/\partial y$.

This now leads us to the following question. Given two functions $M(t,y)$ and $N(t,y)$, does there exist a function $\phi(t,y)$ such that $M(t,y) = \partial\phi/\partial t$ and $N(t,y) = \partial\phi/\partial y$? Unfortunately, the answer to this question is almost always no as the following theorem shows.

Theorem 1. *Let $M(t,y)$ and $N(t,y)$ be continuous and have continuous partial derivatives with respect to t and y in the rectangle R consisting of those points (t,y) with $a < t < b$ and $c < y < d$. There exists a function $\phi(t,y)$ such that $M(t,y) = \partial\phi/\partial t$ and $N(t,y) = \partial\phi/\partial y$ if, and only if,*

$$\partial M/\partial y = \partial N/\partial t$$

in R .

PROOF. Observe that $M(t,y) = \partial\phi/\partial t$ for some function $\phi(t,y)$ if, and only if,

$$\phi(t,y) = \int M(t,y) dt + h(y) \tag{3}$$

where $h(y)$ is an arbitrary function of y . Taking partial derivatives of both sides of (3) with respect to y , we obtain that

$$\frac{\partial\phi}{\partial y} = \int \frac{\partial M(t,y)}{\partial y} dt + h'(y).$$

Hence, $\partial\phi/\partial y$ will be equal to $N(t,y)$ if, and only if,

$$N(t,y) = \int \frac{\partial M(t,y)}{\partial y} dt + h'(y)$$

or

$$h'(y) = N(t,y) - \int \frac{\partial M(t,y)}{\partial y} dt. \tag{4}$$

Now $h'(y)$ is a function of y alone, while the right-hand side of (4) appears to be a function of both t and y . But a function of y alone cannot be equal to a function of both t and y . Thus Equation (4) makes sense only if the right-hand side is a function of y alone, and this is the case if, and only if,

$$\frac{\partial}{\partial t} \left[N(t,y) - \int \frac{\partial M(t,y)}{\partial y} dt \right] = \frac{\partial N}{\partial t} - \frac{\partial M}{\partial y} = 0.$$

1 First-order differential equations

Hence, if $\partial N/\partial t \neq \partial M/\partial y$, then there is no function $\phi(t,y)$ such that $M = \partial\phi/\partial t$, $N = \partial\phi/\partial y$. On the other hand, if $\partial N/\partial t = \partial M/\partial y$ then we can solve for

$$h(y) = \int \left[N(t,y) - \int \frac{\partial M(t,y)}{\partial y} dt \right] dy.$$

Consequently, $M = \partial\phi/\partial t$, and $N = \partial\phi/\partial y$ with

$$\phi(t,y) = \int M(t,y) dt + \int \left[N(t,y) - \int \frac{\partial M(t,y)}{\partial y} dt \right] dy. \quad \square \quad (5)$$

Definition. The differential equation

$$M(t,y) + N(t,y) \frac{dy}{dt} = 0 \quad (6)$$

is said to be *exact* if $\partial M/\partial y = \partial N/\partial t$.

The reason for this definition, of course, is that the left-hand side of (6) is the exact derivative of a known function of t and y if $\partial M/\partial y = \partial N/\partial t$.

Remark 1. It is not essential, in the statement of Theorem 1, that $\partial M/\partial y = \partial N/\partial t$ in a rectangle. It is sufficient if $\partial M/\partial y = \partial N/\partial t$ in any region R which contains no “holes”. That is to say, if C is any closed curve lying entirely in R , then its interior also lies entirely in R .

Remark 2. The differential equation $dy/dt = f(t,y)$ can always be written in the form $M(t,y) + N(t,y)(dy/dt) = 0$ by setting $M(t,y) = -f(t,y)$ and $N(t,y) = 1$.

Remark 3. It is customary to say that the solution of an exact differential equation is given by $\phi(t,y) = \text{constant}$. What we really mean is that the equation $\phi(t,y) = c$ is to be solved for y as a function of t and c . Unfortunately, most exact differential equations cannot be solved explicitly for y as a function of t . While this may appear to be very disappointing, we wish to point out that it is quite simple, with the aid of a computer, to compute $y(t)$ to any desired accuracy (see Section 1.11).

In practice, we do not recommend memorizing Equation (5). Rather, we will follow one of three different methods to obtain $\phi(t,y)$.

First Method: The equation $M(t,y) = \partial\phi/\partial t$ determines $\phi(t,y)$ up to an arbitrary function of y alone, that is,

$$\phi(t,y) = \int M(t,y) dt + h(y).$$

The function $h(y)$ is then determined from the equation

$$h'(y) = N(t, y) - \int \frac{\partial M(t, y)}{\partial y} dt.$$

Second Method: If $N(t, y) = \partial\phi/\partial y$, then, of necessity,

$$\phi(t, y) = \int N(t, y) dy + k(t)$$

where $k(t)$ is an arbitrary function of t alone. Since

$$M(t, y) = \frac{\partial\phi}{\partial t} = \int \frac{\partial N(t, y)}{\partial t} dy + k'(t)$$

we see that $k(t)$ is determined from the equation

$$k'(t) = M(t, y) - \int \frac{\partial N(t, y)}{\partial t} dy.$$

Note that the right-hand side of this equation (see Exercise 2) is a function of t alone if $\partial M/\partial y = \partial N/\partial t$.

Third Method: The equations $\partial\phi/\partial t = M(t, y)$ and $\partial\phi/\partial y = N(t, y)$ imply that

$$\phi(t, y) = \int M(t, y) dt + h(y) \quad \text{and} \quad \phi(t, y) = \int N(t, y) dy + k(t).$$

Usually, we can determine $h(y)$ and $k(t)$ just by inspection.

Example 3. Find the general solution of the differential equation

$$3y + e^t + (3t + \cos y) \frac{dy}{dt} = 0.$$

Solution. Here $M(t, y) = 3y + e^t$ and $N(t, y) = 3t + \cos y$. This equation is exact since $\partial M/\partial y = 3$ and $\partial N/\partial t = 3$. Hence, there exists a function $\phi(t, y)$ such that

$$(i) \quad 3y + e^t = \frac{\partial\phi}{\partial t} \quad \text{and} \quad (ii) \quad 3t + \cos y = \frac{\partial\phi}{\partial y}.$$

We will find $\phi(t, y)$ by each of the three methods outlined above.

First Method: From (i), $\phi(t, y) = e^t + 3ty + h(y)$. Differentiating this equation with respect to y and using (ii) we obtain that

$$h'(y) + 3t = 3t + \cos y.$$

Thus, $h(y) = \sin y$ and $\phi(t, y) = e^t + 3ty + \sin y$. (Strictly speaking, $h(y) = \sin y + \text{constant}$. However, we already incorporate this constant of integration into the solution when we write $\phi(t, y) = c$.) The general solution of the differential equation must be left in the form $e^t + 3ty + \sin y = c$ since we cannot find y explicitly as a function of t from this equation.

Second Method: From (ii), $\phi(t, y) = 3ty + \sin y + k(t)$. Differentiating this

expression with respect to t , and using (i) we obtain that

$$3y + k'(t) = 3y + e^t.$$

Thus, $k(t) = e^t$ and $\phi(t, y) = 3ty + \sin y + e^t$.

Third Method: From (i) and (ii)

$$\phi(t, y) = e^t + 3ty + h(y) \quad \text{and} \quad \phi(t, y) = 3ty + \sin y + k(t).$$

Comparing these two expressions for the *same* function $\phi(t, y)$ it is obvious that $h(y) = \sin y$ and $k(t) = e^t$. Hence

$$\phi(t, y) = e^t + 3ty + \sin y.$$

Example 4. Find the solution of the initial-value problem

$$3t^2y + 8ty^2 + (t^3 + 8t^2y + 12y^2) \frac{dy}{dt} = 0, \quad y(2) = 1.$$

Solution. Here $M(t, y) = 3t^2y + 8ty^2$ and $N(t, y) = t^3 + 8t^2y + 12y^2$. This equation is exact since

$$\frac{\partial M}{\partial y} = 3t^2 + 16ty \quad \text{and} \quad \frac{\partial N}{\partial t} = 3t^2 + 16ty.$$

Hence, there exists a function $\phi(t, y)$ such that

$$(i) \quad 3t^2y + 8ty^2 = \frac{\partial \phi}{\partial t} \quad \text{and} \quad (ii) \quad t^3 + 8t^2y + 12y^2 = \frac{\partial \phi}{\partial y}.$$

Again, we will find $\phi(t, y)$ by each of three methods.

First Method: From (i), $\phi(t, y) = t^3y + 4t^2y^2 + h(y)$. Differentiating this equation with respect to y and using (ii) we obtain that

$$t^3 + 8t^2y + h'(y) = t^3 + 8t^2y + 12y^2.$$

Hence, $h(y) = 4y^3$ and the general solution of the differential equation is $\phi(t, y) = t^3y + 4t^2y^2 + 4y^3 = c$. Setting $t = 2$ and $y = 1$ in this equation, we see that $c = 28$. Thus, the solution of our initial-value problem is defined implicitly by the equation $t^3y + 4t^2y^2 + 4y^3 = 28$.

Second Method: From (ii), $\phi(t, y) = t^3y + 4t^2y^2 + 4y^3 + k(t)$. Differentiating this expression with respect to t and using (i) we obtain that

$$3t^2y + 8ty^2 + k'(t) = 3t^2y + 8ty^2.$$

Thus $k(t) = 0$ and $\phi(t, y) = t^3y + 4t^2y^2 + 4y^3$.

Third Method: From (i) and (ii)

$$\phi(t, y) = t^3y + 4t^2y^2 + h(y) \quad \text{and} \quad \phi(t, y) = t^3y + 4t^2y^2 + 4y^3 + k(t).$$

Comparing these two expressions for the same function $\phi(t, y)$ we see that $h(y) = 4y^3$ and $k(t) = 0$. Hence, $\phi(t, y) = t^3y + 4t^2y^2 + 4y^3$.

In most instances, as Examples 3 and 4 illustrate, the third method is the simplest to use. However, if it is much easier to integrate N with re-

spect to y than it is to integrate M with respect to t , we should use the second method, and vice-versa.

Example 5. Find the solution of the initial-value problem

$$4t^3e^{t+y} + t^4e^{t+y} + 2t + (t^4e^{t+y} + 2y)\frac{dy}{dt} = 0, \quad y(0) = 1.$$

Solution. This equation is exact since

$$\frac{\partial}{\partial y}(4t^3e^{t+y} + t^4e^{t+y} + 2t) = (t^4 + 4t^3)e^{t+y} = \frac{\partial}{\partial t}(t^4e^{t+y} + 2y).$$

Hence, there exists a function $\phi(t, y)$ such that

$$(i) \quad 4t^3e^{t+y} + t^4e^{t+y} + 2t = \frac{\partial\phi}{\partial t}$$

and

$$(ii) \quad t^4e^{t+y} + 2y = \frac{\partial\phi}{\partial y}.$$

Since it is much simpler to integrate $t^4e^{t+y} + 2y$ with respect to y than it is to integrate $4t^3e^{t+y} + t^4e^{t+y} + 2t$ with respect to t , we use the second method. From (ii), $\phi(t, y) = t^4e^{t+y} + y^2 + k(t)$. Differentiating this expression with respect to t and using (i) we obtain

$$(t^4 + 4t^3)e^{t+y} + k'(t) = 4t^3e^{t+y} + t^4e^{t+y} + 2t.$$

Thus, $k(t) = t^2$ and the general solution of the differential equation is $\phi(t, y) = t^4e^{t+y} + y^2 + t^2 = c$. Setting $t=0$ and $y=1$ in this equation yields $c = 1$. Thus, the solution of our initial-value problem is defined implicitly by the equation $t^4e^{t+y} + t^2 + y^2 = 1$.

Suppose now that we are given a differential equation

$$M(t, y) + N(t, y)\frac{dy}{dt} = 0 \tag{7}$$

which is not exact. Can we make it exact? More precisely, can we find a function $\mu(t, y)$ such that the equivalent differential equation

$$\mu(t, y)M(t, y) + \mu(t, y)N(t, y)\frac{dy}{dt} = 0 \tag{8}$$

is exact? This question is simple, in principle, to answer. The condition that (8) be exact is that

$$\frac{\partial}{\partial y}(\mu(t, y)M(t, y)) = \frac{\partial}{\partial t}(\mu(t, y)N(t, y))$$

or

$$M\frac{\partial\mu}{\partial y} + \mu\frac{\partial M}{\partial y} = N\frac{\partial\mu}{\partial t} + \mu\frac{\partial N}{\partial t}. \tag{9}$$

(For simplicity of writing, we have suppressed the dependence of μ, M and N on t and y in (9).) Thus, Equation (8) is exact if and only if $\mu(t, y)$ satisfies Equation (9).

Definition. A function $\mu(t,y)$ satisfying Equation (9) is called an *integrating factor* for the differential equation (7).

The reason for this definition, of course, is that if μ satisfies (9) then we can write (8) in the form $(d/dt)\phi(t,y)=0$ and this equation can be integrated immediately to yield the solution $\phi(t,y)=c$. Unfortunately, though, there are only two special cases where we can find an explicit solution of (9). These occur when the differential equation (7) has an integrating factor which is either a function of t alone, or a function of y alone. Observe that if μ is a function of t alone, then Equation (9) reduces to

$$N \frac{d\mu}{dt} = \mu \left(\frac{\partial M}{\partial y} - \frac{\partial N}{\partial t} \right) \quad \text{or} \quad \frac{d\mu}{dt} = \frac{\left(\frac{\partial M}{\partial y} - \frac{\partial N}{\partial t} \right)}{N} \mu.$$

But this equation is meaningless unless the expression

$$\frac{\frac{\partial M}{\partial y} - \frac{\partial N}{\partial t}}{N}$$

is a function of t alone, that is,

$$\frac{\frac{\partial M}{\partial y} - \frac{\partial N}{\partial t}}{N} = R(t).$$

If this is the case then $\mu(t) = \exp\left(\int R(t) dt\right)$ is an integrating factor for the differential equation (7).

Remark. It should be noted that the expression

$$\frac{\frac{\partial M}{\partial y} - \frac{\partial N}{\partial t}}{N}$$

is almost always a function of both t and y . Only for very special pairs of functions $M(t,y)$ and $N(t,y)$ is it a function of t alone. A similar situation occurs if μ is a function of y alone (see Exercise 17). It is for this reason that we cannot solve very many differential equations.

Example 6. Find the general solution of the differential equation

$$\frac{y^2}{2} + 2ye^t + (y + e^t) \frac{dy}{dt} = 0.$$

Solution. Here $M(t,y) = (y^2/2) + 2ye^t$ and $N(t,y) = y + e^t$. This equation is

not exact since $\partial M/\partial y = y + 2e^t$ and $\partial N/\partial t = e^t$. However,

$$\frac{1}{N} \left(\frac{\partial M}{\partial y} - \frac{\partial N}{\partial t} \right) = \frac{y + e^t}{y + e^t} = 1.$$

Hence, this equation has $\mu(t) = \exp\left(\int 1 dt\right) = e^t$ as an integrating factor.

This means, of course, that the equivalent differential equation

$$\frac{y^2}{2} e^t + 2ye^{2t} + (ye^t + e^{2t}) \frac{dy}{dt} = 0$$

is exact. Therefore, there exists a function $\phi(t, y)$ such that

$$(i) \quad \frac{y^2}{2} e^t + 2ye^{2t} = \frac{\partial \phi}{\partial t}$$

and

$$(ii) \quad ye^t + e^{2t} = \frac{\partial \phi}{\partial y}.$$

From Equations (i) and (ii),

$$\phi(t, y) = \frac{y^2}{2} e^t + ye^{2t} + h(y)$$

and

$$\phi(t, y) = \frac{y^2}{2} e^t + ye^{2t} + k(t).$$

Thus, $h(y) = 0$, $k(t) = 0$ and the general solution of the differential equation is

$$\phi(t, y) = \frac{y^2}{2} e^t + ye^{2t} = c.$$

Solving this equation for y as a function of t we see that

$$y(t) = -e^t \pm [e^{2t} + 2ce^{-t}]^{1/2}.$$

Example 7. Use the methods of this section to find the general solution of the linear equation $(dy/dt) + a(t)y = b(t)$.

Solution. We write this equation in the form $M(t, y) + N(t, y)(dy/dt) = 0$ with $M(t, y) = a(t)y - b(t)$ and $N(t, y) = 1$. This equation is not exact since $\partial M/\partial y = a(t)$ and $\partial N/\partial t = 0$. However, $((\partial M/\partial y) - (\partial N/\partial t))/N = a(t)$.

Hence, $\mu(t) = \exp\left(\int a(t) dt\right)$ is an integrating factor for the first-order linear equation. Therefore, there exists a function $\phi(t, y)$ such that

$$(i) \quad \mu(t)[a(t)y - b(t)] = \frac{\partial \phi}{\partial t}$$

1 First-order differential equations

and

$$(ii) \mu(t) = \frac{\partial \phi}{\partial y}.$$

Now, observe from (ii) that $\phi(t, y) = \mu(t)y + k(t)$. Differentiating this equation with respect to t and using (i) we see that

$$\mu'(t)y + k'(t) = \mu(t)a(t)y - \mu(t)b(t).$$

But, $\mu'(t) = a(t)\mu(t)$. Consequently, $k'(t) = -\mu(t)b(t)$ and

$$\phi(t, y) = \mu(t)y - \int \mu(t)b(t) dt.$$

Hence, the general solution of the first-order linear equation is

$$\mu(t)y - \int \mu(t)b(t) dt = c,$$

and this is the result we obtained in Section 1.2.

EXERCISES

1. Use the theorem of equality of mixed partial derivatives to show that $\partial M / \partial y = \partial N / \partial t$ if the equation $M(t, y) + N(t, y)(dy/dt) = 0$ is exact.
2. Show that the expression $M(t, y) - \int (\partial N(t, y) / \partial t) dy$ is a function of t alone if $\partial M / \partial y = \partial N / \partial t$.

In each of Problems 3–6 find the general solution of the given differential equation.

3. $2t \sin y + y^3 e^t + (t^2 \cos y + 3y^2 e^t) \frac{dy}{dt} = 0$

4. $1 + (1 + ty)e^{ty} + (1 + t^2 e^{ty}) \frac{dy}{dt} = 0$

5. $y \sec^2 t + \sec t \tan t + (2y + \tan t) \frac{dy}{dt} = 0$

6. $\frac{y^2}{2} - 2ye^t + (y - e^t) \frac{dy}{dt} = 0$

In each of Problems 7–11, solve the given initial-value problem.

7. $2ty^3 + 3t^2y^2 \frac{dy}{dt} = 0, \quad y(1) = 1$

8. $2t \cos y + 3t^2y + (t^3 - t^2 \sin y - y) \frac{dy}{dt} = 0, \quad y(0) = 2$

9. $3t^2 + 4ty + (2y + 2t^2) \frac{dy}{dt} = 0, \quad y(0) = 1$

10. $y(\cos 2t)e^{ty} - 2(\sin 2t)e^{ty} + 2t + (t(\cos 2t)e^{ty} - 3) \frac{dy}{dt} = 0, \quad y(0) = 0$

$$11. \quad 3ty + y^2 + (t^2 + ty) \frac{dy}{dt} = 0, \quad y(2) = 1$$

In each of Problems 12–14, determine the constant a so that the equation is exact, and then solve the resulting equation.

$$12. \quad t + ye^{2ty} + ate^{2ty} \frac{dy}{dt} = 0$$

$$13. \quad \frac{1}{t^2} + \frac{1}{y^2} + \frac{(at+1)}{y^3} \frac{dy}{dt} = 0$$

$$14. \quad e^{at+y} + 3t^2y^2 + (2yt^3 + e^{at+y}) \frac{dy}{dt} = 0$$

15. Show that every separable equation of the form $M(t) + N(y)dy/dt = 0$ is exact.

16. Find all functions $f(t)$ such that the differential equation

$$y^2 \sin t + yf(t)(dy/dt) = 0$$

is exact. Solve the differential equation for these $f(t)$.

17. Show that if $((\partial N/\partial t) - (\partial M/\partial y))/M = Q(y)$, then the differential equation $M(t,y) + N(t,y)dy/dt = 0$ has an integrating factor $\mu(y) = \exp\left(\int Q(y) dy\right)$.

18. The differential equation $f(t)(dy/dt) + t^2 + y = 0$ is known to have an integrating factor $\mu(t) = t$. Find all possible functions $f(t)$.

19. The differential equation $e^t \sec y - \tan y + (dy/dt) = 0$ has an integrating factor of the form $e^{-at} \cos y$ for some constant a . Find a , and then solve the differential equation.

20. The Bernoulli differential equation is $(dy/dt) + a(t)y = b(t)y^n$. Multiplying through by $\mu(t) = \exp\left(\int a(t) dt\right)$, we can rewrite this equation in the form $d/dt(\mu(t)y) = b(t)\mu(t)y^n$. Find the general solution of this equation by finding an appropriate integrating factor. *Hint:* Divide both sides of the equation by an appropriate function of y .

1.10 The existence–uniqueness theorem; Picard iteration

Consider the initial-value problem

$$\frac{dy}{dt} = f(t,y), \quad y(t_0) = y_0 \quad (1)$$

where f is a given function of t and y . Chances are, as the remarks in Section 1.9 indicate, that we will be unable to solve (1) explicitly. This leads us to ask the following questions.

1. How are we to know that the initial-value problem (1) actually has a solution if we can't exhibit it?

1 First-order differential equations

2. How do we know that there is only one solution $y(t)$ of (1)? Perhaps there are two, three, or even infinitely many solutions.
3. Why bother asking the first two questions? After all, what's the use of determining whether (1) has a unique solution if we won't be able to explicitly exhibit it?

The answer to the third question lies in the observation that it is never necessary, in applications, to find the solution $y(t)$ of (1) to more than a finite number of decimal places. Usually, it is more than sufficient to find $y(t)$ to four decimal places. As we shall see in Sections 1.13–17, this can be done quite easily with the aid of a digital computer. In fact, we will be able to compute $y(t)$ to eight, and even sixteen, decimal places. Thus, the knowledge that (1) has a unique solution $y(t)$ is our hunting license to go looking for it.

To resolve the first question, we must establish the existence of a function $y(t)$ whose value at $t = t_0$ is y_0 , and whose derivative at any time t equals $f(t, y(t))$. In order to accomplish this, we must find a theorem which enables us to establish the existence of a function having certain properties, without our having to exhibit this function explicitly. If we search through the Calculus, we find that we encounter such a situation exactly once, and this is in connection with the theory of limits. As we show in Appendix B, it is often possible to prove that a sequence of functions $y_n(t)$ has a limit $y(t)$, without our having to exhibit $y(t)$. For example, we can prove that the sequence of functions

$$y_n(t) = \frac{\sin \pi t}{1^2} + \frac{\sin 2\pi t}{2^2} + \dots + \frac{\sin n\pi t}{n^2}$$

has a limit $y(t)$ even though we cannot exhibit $y(t)$ explicitly. This suggests the following algorithm for proving the existence of a solution $y(t)$ of (1).

- (a) Construct a sequence of functions $y_n(t)$ which come closer and closer to solving (1).
- (b) Show that the sequence of functions $y_n(t)$ has a limit $y(t)$ on a suitable interval $t_0 \leq t \leq t_0 + \alpha$.
- (c) Prove that $y(t)$ is a solution of (1) on this interval.

We now show how to implement this algorithm.

(a) *Construction of the approximating sequence $y_n(t)$*

The problem of finding a sequence of functions that come closer and closer to satisfying a certain equation is one that arises quite often in mathematics. Experience has shown that it is often easiest to resolve this problem when our equation can be written in the special form

$$y(t) = L(t, y(t)), \tag{2}$$

where L may depend explicitly on y , and on integrals of functions of y .

For example, we may wish to find a function $y(t)$ satisfying

$$y(t) = 1 + \sin[t + y(t)],$$

or

$$y(t) = 1 + y^2(t) + \int_0^t y^3(s) ds.$$

In these two cases, $L(t, y(t))$ is an abbreviation for

$$1 + \sin[t + y(t)]$$

and

$$1 + y^2(t) + \int_0^t y^3(s) ds,$$

respectively.

The key to understanding what is special about Equation (2) is to view $L(t, y(t))$ as a “machine” that takes in one function and gives back another one. For example, let

$$L(t, y(t)) = 1 + y^2(t) + \int_0^t y^3(s) ds.$$

If we plug the function $y(t) = t$ into this machine, (that is, if we compute $1 + t^2 + \int_0^t s^3 ds$) then the machine returns to us the function $1 + t^2 + t^4/4$. If we plug the function $y(t) = \cos t$ into this machine, then it returns to us the function

$$1 + \cos^2 t + \int_0^t \cos^3 s ds = 1 + \cos^2 t + \sin t - \frac{\sin^3 t}{3}.$$

According to this viewpoint, we can characterize all solutions $y(t)$ of (2) as those functions $y(t)$ which the machine L leaves unchanged. In other words, if we plug a function $y(t)$ into the machine L , and the machine returns to us this same function, then $y(t)$ is a solution of (2).

We can put the initial-value problem (1) into the special form (2) by integrating both sides of the differential equation $y' = f(t, y)$ with respect to t . Specifically, if $y(t)$ satisfies (1), then

$$\int_{t_0}^t \frac{dy(s)}{ds} ds = \int_{t_0}^t f(s, y(s)) ds$$

so that

$$y(t) = y_0 + \int_{t_0}^t f(s, y(s)) ds. \quad (3)$$

Conversely, if $y(t)$ is continuous and satisfies (3), then $dy/dt = f(t, y(t))$. Moreover, $y(t_0)$ is obviously y_0 . Therefore, $y(t)$ is a solution of (1) if, and only if, it is a continuous solution of (3).

1 First-order differential equations

Equation (3) is called an integral equation, and it is in the special form (2) if we set

$$L(t, y(t)) = y_0 + \int_{t_0}^t f(s, y(s)) ds.$$

This suggests the following scheme for constructing a sequence of “approximate solutions” $y_n(t)$ of (3). Let us start by guessing a solution $y_0(t)$ of (3). The simplest possible guess is $y_0(t) = y_0$. To check whether $y_0(t)$ is a solution of (3), we compute

$$y_1(t) = y_0 + \int_{t_0}^t f(s, y_0(s)) ds.$$

If $y_1(t) = y_0$, then $y(t) = y_0$ is indeed a solution of (3). If not, then we try $y_1(t)$ as our next guess. To check whether $y_1(t)$ is a solution of (3), we compute

$$y_2(t) = y_0 + \int_{t_0}^t f(s, y_1(s)) ds,$$

and so on. In this manner, we define a sequence of functions $y_1(t)$, $y_2(t)$, ..., where

$$y_{n+1}(t) = y_0 + \int_{t_0}^t f(s, y_n(s)) ds. \quad (4)$$

These functions $y_n(t)$ are called successive approximations, or Picard iterates, after the French mathematician Picard who first discovered them. Remarkably, these Picard iterates always converge, on a suitable interval, to a solution $y(t)$ of (3).

Example 1. Compute the Picard iterates for the initial-value problem

$$y' = y, \quad y(0) = 1,$$

and show that they converge to the solution $y(t) = e^t$.

Solution. The integral equation corresponding to this initial-value problem is

$$y(t) = 1 + \int_0^t y(s) ds.$$

Hence, $y_0(t) = 1$

$$y_1(t) = 1 + \int_0^t 1 ds = 1 + t$$

$$y_2(t) = 1 + \int_0^t y_1(s) ds = 1 + \int_0^t (1 + s) ds = 1 + t + \frac{t^2}{2!}$$

and, in general,

$$\begin{aligned} y_n(t) &= 1 + \int_0^t y_{n-1}(s) ds = 1 + \int_0^t \left[1 + s + \dots + \frac{s^{n-1}}{(n-1)!} \right] ds \\ &= 1 + t + \frac{t^2}{2!} + \dots + \frac{t^n}{n!}. \end{aligned}$$

Since $e^t = 1 + t + t^2/2! + \dots$, we see that the Picard iterates $y_n(t)$ converge to the solution $y(t)$ of this initial-value problem.

Example 2. Compute the Picard iterates $y_1(t), y_2(t)$ for the initial-value problem $y' = 1 + y^3$, $y(1) = 1$.

Solution. The integral equation corresponding to this initial-value problem is

$$y(t) = 1 + \int_1^t [1 + y^3(s)] ds.$$

Hence, $y_0(t) = 1$

$$y_1(t) = 1 + \int_1^t (1 + 1) ds = 1 + 2(t - 1)$$

and

$$\begin{aligned} y_2(t) &= 1 + \int_1^t \left\{ 1 + [1 + 2(s - 1)]^3 \right\} ds \\ &= 1 + 2(t - 1) + 3(t - 1)^2 + 4(t - 1)^3 + 2(t - 1)^4. \end{aligned}$$

Notice that it is already quite cumbersome to compute $y_3(t)$.

(b) *Convergence of the Picard iterates*

As was mentioned in Section 1.4, the solutions of nonlinear differential equations may not exist for all time t . Therefore, we cannot expect the Picard iterates $y_n(t)$ of (3) to converge for all t . To provide us with a clue, or estimate, of where the Picard iterates converge, we try to find an interval in which all the $y_n(t)$ are uniformly bounded (that is, $|y_n(t)| \leq K$ for some fixed constant K). Equivalently, we seek a rectangle R which contains the graphs of all the Picard iterates $y_n(t)$. Lemma 1 shows us how to find such a rectangle.

Lemma 1. Choose any two positive numbers a and b , and let R be the rectangle: $t_0 \leq t \leq t_0 + a$, $|y - y_0| \leq b$. Compute

$$M = \max_{(t,y) \text{ in } R} |f(t,y)|, \quad \text{and set} \quad \alpha = \min\left(a, \frac{b}{M}\right).$$

Then,

$$|y_n(t) - y_0| \leq M(t - t_0) \tag{5}$$

for $t_0 \leq t \leq t_0 + \alpha$.

Lemma 1 states that the graph of $y_n(t)$ is sandwiched between the lines $y = y_0 + M(t - t_0)$ and $y = y_0 - M(t - t_0)$, for $t_0 \leq t \leq t_0 + \alpha$. These lines leave the rectangle R at $t = t_0 + a$ if $a \leq b/M$, and at $t = t_0 + b/M$ if $b/M < a$ (see Figures 1a and 1b). In either case, therefore, the graph of $y_n(t)$ is contained in R for $t_0 \leq t \leq t_0 + \alpha$.

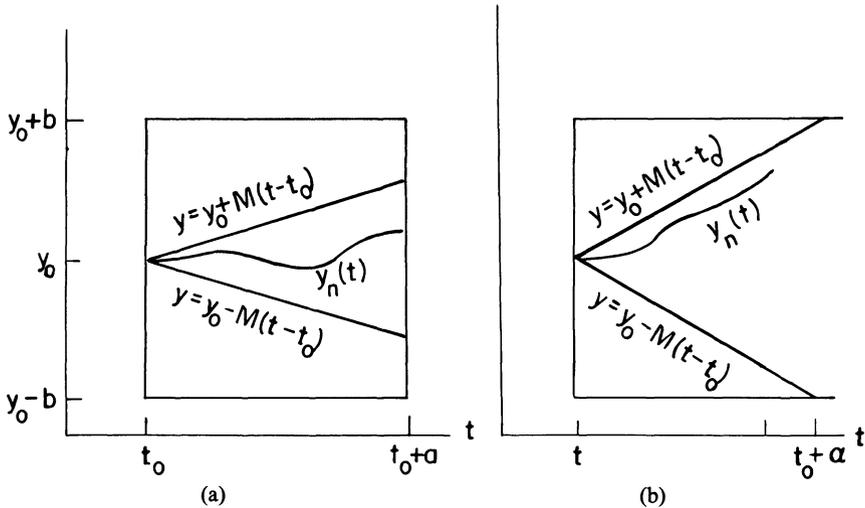


Figure 1. (a) $\alpha = a$; (b) $\alpha = b/M$

PROOF OF LEMMA 1. We establish (5) by induction on n . Observe first that (5) is obviously true for $n=0$, since $y_0(t) = y_0$. Next, we must show that (5) is true for $n=j+1$ if it is true for $n=j$. But this follows immediately, for if $|y_j(t) - y_0| \leq M(t - t_0)$, then

$$\begin{aligned}
 |y_{j+1}(t) - y_0| &= \left| \int_{t_0}^t f(s, y_j(s)) ds \right| \\
 &\leq \int_{t_0}^t |f(s, y_j(s))| ds \leq M(t - t_0)
 \end{aligned}$$

for $t_0 \leq t \leq t_0 + \alpha$. Consequently, (5) is true for all n , by induction. \square

We now show that the Picard iterates $y_n(t)$ of (3) converge for each t in the interval $t_0 \leq t \leq t_0 + \alpha$, if $\partial f/\partial y$ exists and is continuous. Our first step is to reduce the problem of showing that the sequence of functions $y_n(t)$ converges to the much simpler problem of proving that an infinite series converges. This is accomplished by writing $y_n(t)$ in the form

$$y_n(t) = y_0(t) + [y_1(t) - y_0(t)] + \dots + [y_n(t) - y_{n-1}(t)].$$

Clearly, the sequence $y_n(t)$ converges if, and only if, the infinite series

$$[y_1(t) - y_0(t)] + [y_2(t) - y_1(t)] + \dots + [y_n(t) - y_{n-1}(t)] + \dots \quad (6)$$

converges. To prove that the infinite series (6) converges, it suffices to

show that

$$\sum_{n=1}^{\infty} |y_n(t) - y_{n-1}(t)| < \infty. \quad (7)$$

This is accomplished in the following manner. Observe that

$$\begin{aligned} |y_n(t) - y_{n-1}(t)| &= \left| \int_{t_0}^t [f(s, y_{n-1}(s)) - f(s, y_{n-2}(s))] ds \right| \\ &\leq \int_{t_0}^t |f(s, y_{n-1}(s)) - f(s, y_{n-2}(s))| ds \\ &= \int_{t_0}^t \left| \frac{\partial f(s, \xi(s))}{\partial y} \right| |y_{n-1}(s) - y_{n-2}(s)| ds, \end{aligned}$$

where $\xi(s)$ lies between $y_{n-1}(s)$ and $y_{n-2}(s)$. (Recall that $f(x_1) - f(x_2) = f'(\xi)(x_1 - x_2)$, where ξ is some number between x_1 and x_2 .) It follows immediately from Lemma 1 that the points $(s, \xi(s))$ all lie in the rectangle R for $s < t_0 + \alpha$. Consequently,

$$|y_n(t) - y_{n-1}(t)| \leq L \int_{t_0}^t |y_{n-1}(s) - y_{n-2}(s)| ds, \quad t_0 \leq t \leq t_0 + \alpha, \quad (8)$$

where

$$L = \max_{(t,y) \text{ in } R} \left| \frac{\partial f(t,y)}{\partial y} \right|. \quad (9)$$

Equation (9) defines the constant L . Setting $n=2$ in (8) gives

$$\begin{aligned} |y_2(t) - y_1(t)| &\leq L \int_{t_0}^t |y_1(s) - y_0| ds \leq L \int_{t_0}^t M (s - t_0) ds \\ &= \frac{LM(t - t_0)^2}{2}. \end{aligned}$$

This, in turn, implies that

$$\begin{aligned} |y_3(t) - y_2(t)| &\leq L \int_{t_0}^t |y_2(s) - y_1(s)| ds \leq ML^2 \int_{t_0}^t \frac{(s - t_0)^2}{2} ds \\ &= \frac{ML^2(t - t_0)^3}{3!}. \end{aligned}$$

Proceeding inductively, we see that

$$|y_n(t) - y_{n-1}(t)| \leq \frac{ML^{n-1}(t - t_0)^n}{n!}, \quad \text{for } t_0 \leq t \leq t_0 + \alpha. \quad (10)$$

1 First-order differential equations

Therefore, for $t_0 \leq t \leq t_0 + \alpha$,

$$\begin{aligned}
 & |y_1(t) - y_0(t)| + |y_2(t) - y_1(t)| + \dots \\
 & \leq M(t - t_0) + \frac{ML(t - t_0)^2}{2!} + \frac{ML^2(t - t_0)^3}{3!} + \dots \\
 & \leq M\alpha + \frac{ML\alpha^2}{2!} + \frac{ML^2\alpha^3}{3!} + \dots \\
 & = \frac{M}{L} \left[\alpha L + \frac{(\alpha L)^2}{2!} + \frac{(\alpha L)^3}{3!} + \dots \right] \\
 & = \frac{M}{L} (e^{\alpha L} - 1).
 \end{aligned}$$

This quantity, obviously, is less than infinity. Consequently, the Picard iterates $y_n(t)$ converge for each t in the interval $t_0 \leq t \leq t_0 + \alpha$. (A similar argument shows that $y_n(t)$ converges for each t in the interval $t_0 - \beta \leq t \leq t_0$, where $\beta = \min(a, b/N)$, and N is the maximum value of $|f(t, y)|$ for (t, y) in the rectangle $t_0 - a \leq t \leq t_0, |y - y_0| \leq b$.) We will denote the limit of the sequence $y_n(t)$ by $y(t)$. \square

(c) *Proof that $y(t)$ satisfies the initial-value problem (1)*

We will show that $y(t)$ satisfies the integral equation

$$y(t) = y_0 + \int_{t_0}^t f(s, y(s)) ds \tag{11}$$

and that $y(t)$ is continuous. To this end, recall that the Picard iterates $y_n(t)$ are defined recursively through the equation

$$y_{n+1}(t) = y_0 + \int_{t_0}^t f(s, y_n(s)) ds. \tag{12}$$

Taking limits of both sides of (12) gives

$$y(t) = y_0 + \lim_{n \rightarrow \infty} \int_{t_0}^t f(s, y_n(s)) ds. \tag{13}$$

To show that the right-hand side of (13) equals

$$y_0 + \int_{t_0}^t f(s, y(s)) ds,$$

(that is, to justify passing the limit through the integral sign) we must show that

$$\left| \int_{t_0}^t f(s, y(s)) ds - \int_{t_0}^t f(s, y_n(s)) ds \right|$$

approaches zero as n approaches infinity. This is accomplished in the following manner. Observe first that the graph of $y(t)$ lies in the rectangle R for $t \leq t_0 + \alpha$, since it is the limit of functions $y_n(t)$ whose graphs lie in R .

Hence

$$\left| \int_{t_0}^t f(s, y(s)) ds - \int_{t_0}^t f(s, y_n(s)) ds \right| \leq \int_{t_0}^t |f(s, y(s)) - f(s, y_n(s))| ds \leq L \int_{t_0}^t |y(s) - y_n(s)| ds$$

where L is defined by Equation (9). Next, observe that

$$y(s) - y_n(s) = \sum_{j=n+1}^{\infty} [y_j(s) - y_{j-1}(s)]$$

since

$$y(s) = y_0 + \sum_{j=1}^{\infty} [y_j(s) - y_{j-1}(s)]$$

and

$$y_n(s) = y_0 + \sum_{j=1}^n [y_j(s) - y_{j-1}(s)].$$

Consequently, from (10),

$$\begin{aligned} |y(s) - y_n(s)| &\leq M \sum_{j=n+1}^{\infty} L^{j-1} \frac{(s-t_0)^j}{j!} \\ &\leq M \sum_{j=n+1}^{\infty} \frac{L^{j-1} \alpha^j}{j!} = \frac{M}{L} \sum_{j=n+1}^{\infty} \frac{(\alpha L)^j}{j!}, \end{aligned} \quad (14)$$

and

$$\begin{aligned} \left| \int_{t_0}^t f(s, y(s)) ds - \int_{t_0}^t f(s, y_n(s)) ds \right| &\leq M \sum_{j=n+1}^{\infty} \frac{(\alpha L)^j}{j!} \int_{t_0}^t ds \\ &\leq M \alpha \sum_{j=n+1}^{\infty} \frac{(\alpha L)^j}{j!}. \end{aligned}$$

This summation approaches zero as n approaches infinity, since it is the tail end of the convergent Taylor series expansion of $e^{\alpha L}$. Hence,

$$\lim_{n \rightarrow \infty} \int_{t_0}^t f(s, y_n(s)) ds = \int_{t_0}^t f(s, y(s)) ds,$$

and $y(t)$ satisfies (11).

To show that $y(t)$ is continuous, we must show that for every $\varepsilon > 0$ we can find $\delta > 0$ such that

$$|y(t+h) - y(t)| < \varepsilon \quad \text{if } |h| < \delta.$$

Now, we cannot compare $y(t+h)$ with $y(t)$ directly, since we do not know $y(t)$ explicitly. To overcome this difficulty, we choose a large integer N and

1 First-order differential equations

observe that

$$y(t+h) - y(t) = [y(t+h) - y_N(t+h)] \\ + [y_N(t+h) - y_N(t)] + [y_N(t) - y(t)].$$

Specifically, we choose N so large that

$$\frac{M}{L} \sum_{j=N+1}^{\infty} \frac{(\alpha L)^j}{j!} < \frac{\varepsilon}{3}.$$

Then, from (14),

$$|y(t+h) - y_N(t+h)| < \frac{\varepsilon}{3} \quad \text{and} \quad |y_N(t) - y(t)| < \frac{\varepsilon}{3},$$

for $t < t_0 + \alpha$, and h sufficiently small (so that $t+h < t_0 + \alpha$.) Next, observe that $y_N(t)$ is continuous, since it is obtained from N repeated integrations of continuous functions. Therefore, we can choose $\delta > 0$ so small that

$$|y_N(t+h) - y_N(t)| < \frac{\varepsilon}{3} \quad \text{for } |h| < \delta.$$

Consequently,

$$|y(t+h) - y(t)| \leq |y(t+h) - y_N(t+h)| + |y_N(t+h) - y_N(t)| \\ + |y_N(t) - y(t)| < \frac{\varepsilon}{3} + \frac{\varepsilon}{3} + \frac{\varepsilon}{3} = \varepsilon$$

for $|h| < \delta$. Therefore, $y(t)$ is a continuous solution of the integral equation (11), and this completes our proof that $y(t)$ satisfies (1). \square

In summary, we have proven the following theorem.

Theorem 2. *Let f and $\partial f / \partial y$ be continuous in the rectangle $R: t_0 \leq t \leq t_0 + a$, $|y - y_0| \leq b$. Compute*

$$M = \max_{(t,y) \text{ in } R} |f(t,y)|, \quad \text{and set} \quad \alpha = \min\left(a, \frac{b}{M}\right).$$

Then, the initial-value problem $y' = f(t,y)$, $y(t_0) = y_0$ has at least one solution $y(t)$ on the interval $t_0 \leq t \leq t_0 + \alpha$. A similar result is true for $t < t_0$.

Remark. The number α in Theorem 2 depends specifically on our choice of a and b . Different choices of a and b lead to different values of α . Moreover, α doesn't necessarily increase when a and b increase, since an increase in a or b will generally result in an increase in M .

Finally, we turn our attention to the problem of uniqueness of solutions of (1). Consider the initial-value problem

$$\frac{dy}{dt} = (\sin 2t)y^{1/3}, \quad y(0) = 0. \quad (15)$$

One solution of (15) is $y(t) = 0$. Additional solutions can be obtained if we

ignore the fact that $y(0)=0$ and rewrite the differential equation in the form

$$\frac{1}{y^{1/3}} \frac{dy}{dt} = \sin 2t,$$

or

$$\frac{d}{dt} \frac{3y^{2/3}}{2} = \sin 2t.$$

Then,

$$\frac{3y^{2/3}}{2} = \frac{1 - \cos 2t}{2} = \sin^2 t$$

and $y = \pm \sqrt{8/27} \sin^3 t$ are two additional solutions of (15).

Now, initial-value problems that have more than one solution are clearly unacceptable in applications. Therefore, it is important for us to find out exactly what is “wrong” with the initial-value problem (15) that it has more than one solution. If we look carefully at the right-hand side of this differential equation, we see that it does not have a partial derivative with respect to y at $y=0$. This is indeed the problem, as the following theorem shows.

Theorem 2’. *Let f and $\partial f/\partial y$ be continuous in the rectangle $R: t_0 \leq t \leq t_0 + a$, $|y - y_0| \leq b$. Compute*

$$M = \max_{(t,y) \text{ in } R} |f(t,y)|, \quad \text{and set } \alpha = \min\left(a, \frac{b}{M}\right).$$

Then, the initial-value problem

$$y' = f(t,y), \quad y(t_0) = y_0 \tag{16}$$

has a unique solution $y(t)$ on the interval $t_0 \leq t \leq t_0 + \alpha$. In other words, if $y(t)$ and $z(t)$ are two solutions of (16), then $y(t)$ must equal $z(t)$ for $t_0 \leq t \leq t_0 + \alpha$.

PROOF. Theorem 2 guarantees the existence of at least one solution $y(t)$ of (16). Suppose that $z(t)$ is a second solution of (16). Then,

$$y(t) = y_0 + \int_{t_0}^t f(s, y(s)) ds \quad \text{and} \quad z(t) = y_0 + \int_{t_0}^t f(s, z(s)) ds.$$

Subtracting these two equations gives

$$\begin{aligned} |y(t) - z(t)| &= \left| \int_{t_0}^t [f(s, y(s)) - f(s, z(s))] ds \right| \\ &\leq \int_{t_0}^t |f(s, y(s)) - f(s, z(s))| ds \\ &\leq L \int_{t_0}^t |y(s) - z(s)| ds \end{aligned}$$

where L is the maximum value of $|\partial f/\partial y|$ for (t, y) in R . As Lemma 2 below shows, this inequality implies that $y(t) = z(t)$. Hence, the initial-value problem (16) has a unique solution $y(t)$. \square

Lemma 2. *Let $w(t)$ be a nonnegative function, with*

$$w(t) \leq L \int_{t_0}^t w(s) ds. \quad (17)$$

Then, $w(t)$ is identically zero.

FAKE PROOF. Differentiating both sides of (17) gives

$$\frac{dw}{dt} \leq Lw(t), \quad \text{or} \quad \frac{dw}{dt} - Lw(t) \leq 0.$$

Multiplying both sides of this inequality by the integrating factor $e^{-L(t-t_0)}$ gives

$$\frac{d}{dt} e^{-L(t-t_0)} w(t) \leq 0, \quad \text{so that} \quad e^{-L(t-t_0)} w(t) \leq w(t_0)$$

for $t \geq t_0$. But $w(t_0)$ must be zero if $w(t)$ is nonnegative and satisfies (17). Consequently, $e^{-L(t-t_0)} w(t) \leq 0$, and this implies that $w(t)$ is identically zero.

The error in this proof, of course, is that we cannot differentiate both sides of an inequality, and still expect to preserve the inequality. For example, the function $f_1(t) = 2t - 2$ is less than $f_2(t) = t$ on the interval $[0, 1]$, but $f_1'(t)$ is greater than $f_2'(t)$ on this interval. We make this proof “kosher” by the clever trick of setting

$$U(t) = \int_{t_0}^t w(s) ds.$$

Then,

$$\frac{dU}{dt} = w(t) \leq L \int_{t_0}^t w(s) ds = LU(t).$$

Consequently, $e^{-L(t-t_0)} U(t) \leq U(t_0) = 0$, for $t \geq t_0$, and thus $U(t) = 0$. This, in turn, implies that $w(t) = 0$ since

$$0 \leq w(t) \leq L \int_{t_0}^t w(s) ds = LU(t) = 0. \quad \square$$

Example 3. Show that the solution $y(t)$ of the initial-value problem

$$\frac{dy}{dt} = t^2 + e^{-y^2}, \quad y(0) = 0$$

exists for $0 \leq t \leq \frac{1}{2}$, and in this interval, $|y(t)| \leq 1$.

Solution. Let R be the rectangle $0 \leq t \leq \frac{1}{2}$, $|y| \leq 1$. Computing

$$M = \max_{(t,y) \text{ in } R} t^2 + e^{-y^2} = 1 + \left(\frac{1}{2}\right)^2 = \frac{5}{4},$$

we see that $y(t)$ exists for

$$0 \leq t \leq \min\left(\frac{1}{2}, \frac{1}{5/4}\right) = \frac{1}{2},$$

and in this interval, $|y(t)| \leq 1$.

Example 4. Show that the solution $y(t)$ of the initial-value problem

$$\frac{dy}{dt} = e^{-t^2} + y^3, \quad y(0) = 1$$

exists for $0 \leq t \leq 1/9$, and in this interval, $0 \leq y \leq 2$.

Solution. Let R be the rectangle $0 \leq t \leq \frac{1}{9}$, $0 \leq y \leq 2$. Computing

$$M = \max_{(t,y) \text{ in } R} e^{-t^2} + y^3 = 1 + 2^3 = 9,$$

we see that $y(t)$ exists for

$$0 \leq t \leq \min\left(\frac{1}{9}, \frac{1}{9}\right)$$

and in this interval, $0 \leq y \leq 2$.

Example 5. What is the largest interval of existence that Theorem 2 predicts for the solution $y(t)$ of the initial-value problem $y' = 1 + y^2$, $y(0) = 0$?

Solution. Let R be the rectangle $0 \leq t \leq a$, $|y| \leq b$. Computing

$$M = \max_{(t,y) \text{ in } R} 1 + y^2 = 1 + b^2,$$

we see that $y(t)$ exists for

$$0 \leq t \leq \alpha = \min\left(a, \frac{b}{1 + b^2}\right).$$

Clearly, the largest α that we can achieve is the maximum value of the function $b/(1 + b^2)$. This maximum value is $\frac{1}{2}$. Hence, Theorem 2 predicts that $y(t)$ exists for $0 \leq t \leq \frac{1}{2}$. The fact that $y(t) = \tan t$ exists for $0 \leq t < \pi/2$ points out the limitation of Theorem 2.

Example 6. Suppose that $|f(t,y)| \leq K$ in the strip $t_0 \leq t < \infty$, $-\infty < y < \infty$. Show that the solution $y(t)$ of the initial-value problem $y' = f(t,y)$, $y(t_0) = y_0$ exists for all $t \geq t_0$.

Solution. Let R be the rectangle $t_0 \leq t \leq t_0 + a$, $|y - y_0| \leq b$. The quantity

$$M = \max_{(t,y) \text{ in } R} |f(t,y)|$$

1 First-order differential equations

is at most K . Hence, $y(t)$ exists for

$$t_0 \leq t \leq t_0 + \min(a, b/K).$$

Now, we can make the quantity $\min(a, b/K)$ as large as desired by choosing a and b sufficiently large. Therefore $y(t)$ exists for $t \geq t_0$.

EXERCISES

- Construct the Picard iterates for the initial-value problem $y' = 2t(y + 1)$, $y(0) = 0$ and show that they converge to the solution $y(t) = e^{t^2} - 1$.
- Compute the first two Picard iterates for the initial-value problem $y' = t^2 + y^2$, $y(0) = 1$.
- Compute the first three Picard iterates for the initial-value problem $y' = e^t + y^2$, $y(0) = 0$.

In each of Problems 4–15, show that the solution $y(t)$ of the given initial-value problem exists on the specified interval.

- $y' = y^2 + \cos t^2$, $y(0) = 0$; $0 \leq t \leq \frac{1}{2}$
- $y' = 1 + y + y^2 \cos t$, $y(0) = 0$; $0 \leq t \leq \frac{1}{3}$
- $y' = t + y^2$, $y(0) = 0$; $0 \leq t \leq (\frac{1}{2})^{2/3}$
- $y' = e^{-t^2} + y^2$, $y(0) = 0$; $0 \leq t \leq \frac{1}{2}$
- $y' = e^{-t^2} + y^2$, $y(1) = 0$; $1 \leq t \leq 1 + \sqrt{e}/2$
- $y' = e^{-t^2} + y^2$, $y(0) = 1$; $0 \leq t \leq \frac{\sqrt{2}}{1 + (1 + \sqrt{2})^2}$
- $y' = y + e^{-y} + e^{-t}$, $y(0) = 0$; $0 \leq t \leq 1$
- $y' = y^3 + e^{-5t}$, $y(0) = 0.4$; $0 \leq t \leq \frac{3}{10}$
- $y' = e^{(y-t)^2}$, $y(0) = 1$; $0 \leq t \leq \frac{\sqrt{3}-1}{2} e^{-((1+\sqrt{3})/2)^2}$
- $y' = (4y + e^{-t^2})e^{2y}$, $y(0) = 0$; $0 \leq t \leq \frac{1}{8\sqrt{e}}$
- $y' = e^{-t} + \ln(1 + y^2)$, $y(0) = 0$; $0 \leq t < \infty$
- $y' = \frac{1}{4}(1 + \cos 4t)y - \frac{1}{800}(1 - \cos 4t)y^2$, $y(0) = 100$; $0 \leq t \leq 1$
- Consider the initial-value problem

$$y' = t^2 + y^2, \quad y(0) = 0, \quad (*)$$

and let R be the rectangle $0 \leq t \leq a$, $-b \leq y \leq b$.

(a) Show that the solution $y(t)$ of (*) exists for

$$0 \leq t \leq \min\left(a, \frac{b}{a^2 + b^2}\right).$$

1.11 Finding roots of equations by iteration

(b) Show that the maximum value of $b/(a^2 + b^2)$, for a fixed, is $1/2a$.

(c) Show that $\alpha = \min(a, \frac{1}{2}a)$ is largest when $a = 1/\sqrt{2}$.

(d) Conclude that the solution $y(t)$ of (*) exists for $0 \leq t \leq 1/\sqrt{2}$.

17. Prove that $y(t) = -1$ is the only solution of the initial-value problem

$$y' = t(1 + y), \quad y(0) = -1.$$

18. Find a nontrivial solution of the initial-value problem $y' = ty^a$, $y(0) = 0$, $a > 1$. Does this violate Theorem 2'? Explain.

19. Find a solution of the initial-value problem $y' = t\sqrt{1 - y^2}$, $y(0) = 1$, other than $y(t) = 1$. Does this violate Theorem 2'? Explain.

20. Here is an alternate proof of Lemma 2. Let $w(t)$ be a nonnegative function with

$$w(t) \leq L \int_{t_0}^t w(s) ds \quad (*)$$

on the interval $t_0 \leq t \leq t_0 + \alpha$. Since $w(t)$ is continuous, we can find a constant A such that $0 \leq w(t) \leq A$ for $t_0 \leq t \leq t_0 + \alpha$.

(a) Show that $w(t) \leq LA(t - t_0)$.

(b) Use this estimate of $w(t)$ in (*) to obtain

$$w(t) \leq \frac{AL^2(t - t_0)^2}{2}.$$

(c) Proceeding inductively, show that $w(t) \leq AL^n(t - t_0)^n/n!$, for every integer n .

(d) Conclude that $w(t) = 0$ for $t_0 \leq t \leq t_0 + \alpha$.

1.11 Finding roots of equations by iteration

Suppose that we are interested in finding the roots of an equation having the special form

$$x = f(x). \quad (1)$$

For example, we might want to find the roots of the equation

$$x = \sin x + \frac{1}{4}.$$

The methods introduced in the previous section suggest the following algorithm for solving this problem.

1. Try an initial guess x_0 , and use this number to construct a sequence of guesses x_1, x_2, x_3, \dots , where $x_1 = f(x_0)$, $x_2 = f(x_1)$, $x_3 = f(x_2)$, and so on.
2. Show that this sequence of iterates x_n has a limit η as n approaches infinity.
3. Show that η is a root of (1); i.e., $\eta = f(\eta)$.

The following theorem tells us when this algorithm will work.

Theorem 3. *Let $f(x)$ and $f'(x)$ be continuous in the interval $a \leq x \leq b$, with $|f'(x)| \leq \lambda < 1$ in this interval. Suppose, moreover, that the iterates x_n , de-*

defined recursively by the equation

$$x_{n+1} = f(x_n) \quad (2)$$

all lie in the interval $[a, b]$. Then, the iterates x_n converge to a unique number η satisfying (1).

PROOF. We convert the problem of proving that the sequence x_n converges to the simpler problem of proving that an infinite series converges by writing x_n in the form

$$x_n = x_0 + (x_1 - x_0) + (x_2 - x_1) + \dots + (x_n - x_{n-1}).$$

Clearly, the sequence x_n converges if, and only if, the infinite series

$$(x_1 - x_0) + (x_2 - x_1) + \dots + (x_n - x_{n-1}) + \dots = \sum_{n=1}^{\infty} (x_n - x_{n-1})$$

converges. To prove that this infinite series converges, it suffices to show that

$$|x_1 - x_0| + |x_2 - x_1| + \dots = \sum_{n=1}^{\infty} |x_n - x_{n-1}| < \infty.$$

This is accomplished in the following manner. By definition, $x_n = f(x_{n-1})$ and $x_{n-1} = f(x_{n-2})$. Subtracting these two equations gives

$$x_n - x_{n-1} = f(x_{n-1}) - f(x_{n-2}) = f'(\xi)(x_{n-1} - x_{n-2}),$$

where ξ is some number between x_{n-1} and x_{n-2} . In particular, ξ is in the interval $[a, b]$. Therefore, $|f'(\xi)| \leq \lambda$, and

$$|x_n - x_{n-1}| \leq \lambda |x_{n-1} - x_{n-2}|. \quad (3)$$

Iterating this inequality $n - 1$ times gives

$$\begin{aligned} |x_n - x_{n-1}| &\leq \lambda |x_{n-1} - x_{n-2}| \\ &\leq \lambda^2 |x_{n-2} - x_{n-3}| \\ &\vdots \\ &\leq \lambda^{n-1} |x_1 - x_0|. \end{aligned}$$

Consequently,

$$\begin{aligned} \sum_{n=1}^{\infty} |x_n - x_{n-1}| &\leq \sum_{n=1}^{\infty} \lambda^{n-1} |x_1 - x_0| \\ &= |x_1 - x_0| [1 + \lambda + \lambda^2 + \dots] = \frac{|x_1 - x_0|}{1 - \lambda}. \end{aligned}$$

This quantity, obviously, is less than infinity. Therefore, the sequence of iterates x_n has a limit η as n approaches infinity. Taking limits of both

sides of (2) gives

$$\eta = \lim_{n \rightarrow \infty} x_{n+1} = \lim_{n \rightarrow \infty} f(x_n) = f(\eta).$$

Hence, η is a root of (1).

Finally, suppose that η is not unique; that is, there exist two solutions η_1 and η_2 of (1) in the interval $[a, b]$. Then,

$$\eta_1 - \eta_2 = f(\eta_1) - f(\eta_2) = f'(\xi)(\eta_1 - \eta_2),$$

where ξ is some number between η_1 and η_2 . This implies that $\eta_1 = \eta_2$ or $f'(\xi) = 1$. But $f'(\xi)$ cannot be one, since ξ is in the interval $[a, b]$. Therefore, $\eta_1 = \eta_2$. \square

Example 1. Show that the sequence of iterates

$$x_0, \quad x_1 = 1 + \frac{1}{2} \arctan x_0, \quad x_2 = 1 + \frac{1}{2} \arctan x_1, \dots$$

converge to a unique number η satisfying

$$\eta = 1 + \frac{1}{2} \arctan \eta$$

for every initial guess x_0 .

Solution. Let $f(x) = 1 + \frac{1}{2} \arctan x$. Computing $f'(x) = \frac{1}{2} \cdot 1/(1+x^2)$, we see that $|f'(x)|$ is always less than or equal to $\frac{1}{2}$. Hence, by Theorem 3, the sequence of iterates x_0, x_1, x_2, \dots converges to the unique root η of the equation $x = 1 + \frac{1}{2} \arctan x$, for every choice of x_0 .

There are many instances where we know, a priori, that the equation $x = f(x)$ has a unique solution η in a given interval $[a, b]$. In these instances, we can use Theorem 3 to obtain a very good approximation of η . Indeed, life is especially simple in these instances, since we don't have to check that the iterates x_n all lie in a specified interval. If x_0 is sufficiently close to η , then the iterates x_n will always converge to η , as we now show.

Theorem 4. Assume that $f(\eta) = \eta$, and that $|f'(x)| \leq \lambda < 1$ in the interval $|x - \eta| \leq \alpha$. Choose a number x_0 in this interval. Then, the sequence of iterates x_n , defined recursively by the equation $x_{n+1} = f(x_n)$, will always converge to η .

PROOF. Denote the interval $|x - \eta| \leq \alpha$ by I . By Theorem 3, it suffices to show that all the iterates x_n lie in I . To this end, observe that

$$x_{j+1} - \eta = f(x_j) - f(\eta) = f'(\xi)(x_j - \eta)$$

where ξ is some number between x_j and η . In particular, ξ is in I if x_j is in

1 First-order differential equations

I. Thus,

$$|x_{j+1} - \eta| \leq \lambda |x_j - \eta| < |x_j - \eta| \quad (4)$$

if x_j is in I . This implies that x_{j+1} is in I whenever x_j is in I . By induction, therefore, all the iterates x_n lie in I . \square

Equation (4) also shows that x_{n+1} is closer to η than x_n . Specifically, the error we make in approximating η by x_n decreases by at least a factor of λ each time we increase n . Thus, if λ is very small, then the convergence of x_n to η is very rapid, while if λ is close to one, then the convergence is very slow.

Example 2.

(a) Show that the equation

$$x = \sin x + \frac{1}{4} \quad (5)$$

has a unique root η in the interval $[\pi/4, \pi/2]$.

(b) Show that the sequence of numbers

$$x_0, \quad x_1 = \sin x_0 + \frac{1}{4}, \quad x_2 = \sin x_1 + \frac{1}{4}, \dots$$

will converge to η if $\pi/4 \leq x_0 \leq \pi/2$.

(c) Write a computer program to evaluate the first N iterates x_1, x_2, \dots, x_N .
Solution.

(a) Let $g(x) = x - \sin x - \frac{1}{4}$, and observe that $g(\pi/4)$ is negative while $g(\pi/2)$ is positive. Moreover, $g(x)$ is a monotonic increasing function of x for $\pi/4 \leq x \leq \pi/2$, since its derivative is strictly positive in this interval. Therefore, Equation (5) has a unique root $x = \eta$ in the interval $\pi/4 < x < \pi/2$.

(b) Let I denote the interval $\eta - \pi/4 \leq x \leq \eta + \pi/4$. The left endpoint of this interval is greater than zero, while the right endpoint is less than $3\pi/4$. Hence, there exists a number λ , with $0 < \lambda < 1$, such that

$$|\cos x| = \left| \frac{d}{dx} \left(\sin x + \frac{1}{4} \right) \right| \leq \lambda$$

for x in I . Clearly, the interval $[\pi/4, \pi/2]$ is contained in I . Therefore, by Theorem 4, the sequence of numbers

$$x_0, \quad x_1 = \sin x_0 + \frac{1}{4}, \quad x_2 = \sin x_1 + \frac{1}{4}, \dots$$

will converge to η for every x_0 in the interval $[\pi/4, \pi/2]$.

(c)

Pascal Program

```

Program Iterate (input, output);
var
  X: array[0..199] of real;
  k, N: integer;
begin
  readln(X[0], N);
  page;
  writeln('N':4, 'X[N]':14);
  for k := 0 to N do
    begin
      writeln(K:4, ' ':4, X[k]:17:9);
      X[k + 1] := 0.25 + sin(X[k]);
    end;
end.

```

Fortran Program

10	DIMENSION X(200)
	READ (5, 10) X0, N
C	FORMAT (F15.8, I5)
	COMPUTE X(1) FIRST
	X(1) = 0.25 + SIN(X0)
	KA = 0
	KB = 1
	WRITE (6, 20) KA, X0, KB, X(1)
20	FORMAT (1H1, 4X, 'N', 10X, 'X' / (1H, 3X, I3, 4X, F15.9))
C	COMPUTE X(2) THRU X(N)
	DO 40 K = 2, N
	X(K) = 0.25 + SIN(X(K - 1))
	WRITE (6, 30) K, X(K)
30	FORMAT (1H, 3X, I3, 4X, F15.9)
40	CONTINUE
	CALL EXIT
	END

See also C Program 1 in Appendix C for a sample C program.

1 First-order differential equations

Table 1

n	x_n	n	x_n
0	1	8	1.17110411
1	1.09147099	9	1.17122962
2	1.13730626	10	1.17122964
3	1.15750531	11	1.17122965
4	1.16580403	12	1.17122965
5	1.16910543	13	1.17122965
6	1.17040121	14	1.17122965
7	1.17090706	15	1.17122965

In many instances, we want to compute a root of the equation $x = f(x)$ to within a certain accuracy. The easiest, and most efficient way of accomplishing this is to instruct the computer to terminate the program at $k = j$ if x_{j+1} agrees with x_j within the prescribed accuracy.

EXERCISES

- Let η be the unique root of Equation (5).
 - Let $x_0 = \pi/4$. Show that 20 iterations are required to find η to 8 significant decimal places.
 - Let $x_0 = \pi/2$. Show that 20 iterations are required to find η to 8 decimal places.
 - Let $x_0 = 3\pi/8$. Show that 16 iterations are required to find η to 8 decimal places.

- Determine suitable values of x_0 so that the iterates x_n , defined by the equation

$$x_{n+1} = x_n - \frac{1}{4}(x_n^2 - 2)$$

will converge to $\sqrt{2}$.

- Choose $x_0 = 1.4$. Show that 14 iterations are required to find $\sqrt{2}$ to 8 significant decimal places. ($\sqrt{2} = 1.41421356$ to 8 significant decimal places.)
- Determine suitable values of x_0 so that the iterates x_n , defined by the equation

$$x_{n+1} = x_n - \frac{1}{10}(x_n^2 - 2)$$

will converge to $\sqrt{2}$.

- Choose $x_0 = 1.4$. Show that 30 iterations are required to find $\sqrt{2}$ to 6 significant decimal places.

- Determine a suitable value of α so that the iterates x_n , defined by the equation

$$x_{n+1} = x_n - \alpha(x_n^2 - 3), \quad x_0 = 1.7$$

will converge to $\sqrt{3}$.

- Find $\sqrt{3}$ to 6 significant decimal places.

5. Let η be the unique root of the equation $x = 1 + \frac{1}{2} \arctan x$. Find η to 5 significant decimal places.
6. (a) Show that the equation $2 - x = (\ln x)/4$ has a unique root $x = \eta$ in the interval $0 < x < \infty$.
 (b) Let
- $$x_{n+1} = 2 - (\ln x_n)/4, \quad n = 0, 1, 2, \dots$$
- Show that $1 \leq x_n \leq 2$ if $1 \leq x_0 \leq 2$.
 (c) Prove that $x_n \rightarrow \eta$ as $n \rightarrow \infty$ if $1 \leq x_0 \leq 2$.
 (d) Compute η to 5 significant decimal places.
7. (a) Show that the equation $x = \cos x$ has a unique root $x = \eta$ in the interval $0 < x \leq 1$.
 (b) Let $x_{n+1} = \cos x_n$, $n = 0, 1, 2, \dots$, with $0 < x_0 < 1$. Show that $0 < x_n < 1$. Conclude, therefore, that $x_n \rightarrow \eta$ as $n \rightarrow \infty$.
 (c) Find η to 5 significant decimal places.

1.11.1 Newton's method

The method of iteration which we used to solve the equation $x = f(x)$ can also be used to solve the equation $g(x) = 0$. To wit, any solution $x = \eta$ of the equation $g(x) = 0$ is also a solution of the equation

$$x = f(x) = x - g(x), \quad (1)$$

and vice-versa. Better yet, any solution $x = \eta$ of the equation $g(x) = 0$ is also a solution of the equation

$$x = f(x) = x - \frac{g(x)}{h(x)} \quad (2)$$

for any function $h(x)$. Of course, $h(x)$ must be unequal to zero for x near η .

Equation (2) has an arbitrary function $h(x)$ in it. Let us try and choose $h(x)$ so that (i) the assumptions of Theorem 4, Section 1.11 are satisfied, and (ii) the iterates

$$x_0, \quad x_1 = x_0 - \frac{g(x_0)}{h(x_0)}, \quad x_2 = x_1 - \frac{g(x_1)}{h(x_1)}, \dots$$

converge as "rapidly as possible" to the desired root η . To this end, we compute

$$f'(x) = \frac{d}{dx} \left[x - \frac{g(x)}{h(x)} \right] = 1 - \frac{g'(x)}{h(x)} + \frac{h'(x)g(x)}{h^2(x)}$$

and observe that

$$f'(\eta) = 1 - \frac{g'(\eta)}{h(\eta)}.$$

1 First-order differential equations

This suggests that we set $h(x) = g'(x)$, since then $f'(\eta) = 0$. Consequently, the iterates x_n , defined recursively by the equation

$$x_{n+1} = x_n - \frac{g(x_n)}{g'(x_n)}, \quad n = 0, 1, 2, \dots \quad (3)$$

will converge to η if the initial guess x_0 is sufficiently close to η . (If $f'(\eta) = 0$, then $|f'(x)| \leq \lambda < 1$ for $|x - \eta|$ sufficiently small.) Indeed, the choice of $h(x) = f'(x)$ is an *optimal* choice of $h(x)$, since the convergence of x_n to η will be extremely rapid. This follows immediately from the fact that the number λ in Equation 4, Section 1.11 can be taken arbitrarily small, as x_n approaches η .

The iteration scheme (3) is known as Newton's method for solving the equation $g(x) = 0$. It can be shown that if $g(\eta) = 0$, and x_0 is sufficiently close to η , then

$$|x_{n+1} - \eta| \leq c|x_n - \eta|^2,$$

for some positive constant c . In other words, the error we make in approximating η by x_{n+1} is proportional to the square of the error we make in approximating η by x_n . This type of convergence is called quadratic convergence, and it implies that the iterates x_n converge extremely rapidly to η . In many instances, only five or six iterations are required to find η to eight or more significant decimal places.

Example 1. Use Newton's method to compute $\sqrt{2}$.

Solution. The square root of two is a solution of the equation

$$g(x) = x^2 - 2 = 0.$$

Hence, Newton's scheme for this problem is

$$\begin{aligned} x_{n+1} &= x_n - \frac{g(x_n)}{g'(x_n)} = x_n - \frac{(x_n^2 - 2)}{2x_n} \\ &= \frac{x_n}{2} + \frac{1}{x_n}, \quad n = 0, 1, 2, \dots \end{aligned} \quad (4)$$

Sample Pascal and Fortran programs to compute the first N iterates of an initial guess x_0 are given below.

Pascal Program

```
Program Newton (input, output);
var
  X: array[0..199] of real;
  k, N: integer;
begin
  readln(X[0], N);
  page;
```

```

writeIn('N':4, 'X[N]':14);
for k := 0 to N do
  begin
    writeIn(K:4, ' ':4, X[k]:17:9);
    X[k + 1] := X[k]/2 + 1/X[k];
  end;
end.

```

Table 1

n	x_n	n	x_n
0	1.4	3	1.41421356
1	1.41428571	4	1.41421356
2	1.41421356	5	1.41421356

Fortran Program

We need only replace the instructions for computing $X(1)$ and $X(K)$ in the Fortran program of Section 1.11 by

$$X(1) = (X0/2) + 1/X0$$

and

$$X(K) = (X(K-1)/2) + 1/X(K-1)$$

We ran these programs for $x_0 = 1.4$ and $N = 5$, and the results are given in Table 1. Notice that Newton's method requires only 2 iterations to find $\sqrt{2}$ to eight significant decimal places.

See also C Program 2 in Appendix C for a sample C program.

Example 2. Use Newton's method to find the impact velocity of the drums in Section 1.7.

Solution. The impact velocity of the drums satisfies the equation

$$g(v) = v + \frac{300cg}{W} + \frac{W-B}{c} \ln \left[\frac{W-B-cv}{W-B} \right] = 0 \quad (5)$$

where

$$c = 0.08, \quad g = 32.2, \quad W = 527.436, \quad \text{and} \quad B = 470.327.$$

Setting $a = (W-B)/c$ and $d = 300cg/W$ puts (5) in the simpler form

$$g(v) = v + d + a \ln(1 - v/a) = 0. \quad (6)$$

Newton's iteration scheme for this problem is

1 First-order differential equations

$$\begin{aligned}v_{n+1} &= v_n - \frac{g(v_n)}{g'(v_n)} = v_n + \frac{(1 - v_n/a)[v_n + d + a \ln(1 - v_n/a)]}{v_n/a} \\ &= v_n + \frac{a - v_n}{v_n} [v_n + d + a \ln(1 - v_n/a)], \quad n=0, 1, 2, \dots\end{aligned}$$

Sample Pascal and Fortran programs to compute the first N iterates of v_0 are given below.

Pascal Program

Program Newton (input, output);

const

```
c = 0.08;
g = 32.2;
W = 527.436;
B = 470.327;
```

var

```
V: array[0..199] of real;
a, d: real;
k, N: integer;
```

begin

```
readln(V[0], N);
a := (W - B)/c;
d := 300 * c * g/W;
page;
writeln('N':4, 'V[N]':14);
for k := 0 to N do
  begin
    writeln(K:4, ' ':4, V[k]:17:9);
    V[k + 1] := V[k] + ((a - V[k])/V[k])
      * (V[k] + d + a * ln(1 - (V[k]/a)));
  end;
```

end.

Fortran Program

Change every X to V , and replace the instructions for $X(1)$ and $X(K)$ in the Fortran program of Section 1.11 by

$$V(1) = V_0 + ((A - V_0)/V_0) * (V_0 + D + A * A \text{ LOG}(1 - (V_0 / A)))$$

and

$$V(K) = V(K-1) + ((A - V(K-1))/V(K-1)) * (V(K-1) + D + A * A \text{ LOG}(1 - (V(K-1) / A)))$$

(Before running these programs, of course, we must instruct the computer to evaluate the constants $a = (W - B)/c$ and $d = 300 \text{ cg}/W$.)

As was shown in Section 1.7, $v_0 = 45.7$ is a very good approximation of v . We set $v_0 = 45.7$ in the above programs, and the iterates v_n converged very rapidly to $v = 45.1$ ft/s. Thus, the drums can indeed break upon impact.

In general, it is not possible to determine, a priori, how many iterations will be required to achieve a certain accuracy. In practice, we usually take N very large, and instruct the computer to terminate the program if one of the iterates agrees with its predecessor to the desired accuracy.

See also C Program 3 in Appendix C for a sample C program.

EXERCISES

1. Show that the iterates x_n defined by (4) will converge to $\sqrt{2}$ if

$$\sqrt{2/3} < x_0 < \sqrt{2} + (\sqrt{2} - \sqrt{2/3}).$$

2. Use Newton's method to find the following numbers to 8 significant decimal places. (a) $\sqrt{3}$, (b) $\sqrt{5}$, (c) $\sqrt{7}$.
3. The number π is a root of the equation

$$\tan \frac{x}{4} - \cot \frac{x}{4} = 0.$$

Use Newton's method to find π to 8 significant decimal places.

Show that each of the following equations has a unique solution in the given interval, and use Newton's method to find it to 5 significant decimal places.

4. $2x - \tan x = 0$; $\pi < x < 3\pi/2$ 5. $\frac{1}{2} - x + \frac{1}{3}\sin x = 0$; $\frac{1}{2} < x < 1$
 6. $\ln x + (x+1)^3 = 0$; $0 < x < 1$ 7. $2\sqrt{x} = \cos \frac{\pi x}{2}$; $0 < x < 1$
 8. $(x-1)^2 - \frac{1}{2}e^x = 0$; $0 < x < 1$ 9. $x - e^{-x^2} = 1$; $0 < x < 2$.

1.12 Difference equations, and how to compute the interest due on your student loans

In Sections 1.13–1.16 we will construct various approximations of the solution of the initial-value problem $dy/dt = f(t, y)$, $y(t_0) = y_0$. In determining how good these approximations are, we will be confronted with the follow-

ing problem: How large can the numbers E_1, \dots, E_N be if

$$E_{n+1} \leq AE_n + B, \quad n=0, 1, \dots, N-1 \quad (1)$$

for some positive constants A and B , and $E_0=0$? This is a very difficult problem since it deals with *inequalities*, rather than *equalities*. Fortunately, though, we can convert the problem of solving the inequalities (1) into the simpler problem of solving a system of equalities. This is the content of the following lemma.

Lemma 1. *Let E_1, \dots, E_N satisfy the inequalities*

$$E_{n+1} \leq AE_n + B, \quad E_0 = 0$$

for some positive constants A and B . Then, E_n is less than or equal to y_n , where

$$y_{n+1} = Ay_n + B, \quad y_0 = 0. \quad (2)$$

PROOF. We prove Lemma 1 by induction on n . To this end, observe that Lemma 1 is obviously true for $n=0$. Next, we assume that Lemma 1 is true for $n=j$. We must show that Lemma 1 is also true for $n=j+1$. That is to say, we must prove that $E_j \leq y_j$ implies $E_{j+1} \leq y_{j+1}$. But this follows immediately, for if $E_j \leq y_j$ then

$$E_{j+1} \leq AE_j + B \leq Ay_j + B = y_{j+1}.$$

By induction, therefore, $E_n \leq y_n$, $n=0, 1, \dots, N$. □

Our next task is to solve Equation (2), which is often referred to as a difference equation. We will accomplish this in two steps. First we will solve the “simple” difference equation

$$y_{n+1} = y_n + B_n, \quad y_0 = y_0. \quad (3)$$

Then we will reduce the difference equation (2) to the difference equation (3) by a clever change of variables.

Equation (3) is trivial to solve. Observe that

$$\begin{aligned} y_1 - y_0 &= B_0 \\ y_2 - y_1 &= B_1 \\ &\vdots \\ y_{n-1} - y_{n-2} &= B_{n-2} \\ y_n - y_{n-1} &= B_{n-1}. \end{aligned}$$

Adding these equations gives

$$(y_n - y_{n-1}) + (y_{n-1} - y_{n-2}) + \dots + (y_1 - y_0) = B_0 + B_1 + \dots + B_{n-1}.$$

Hence,

$$y_n = y_0 + B_0 + \dots + B_{n-1} = y_0 + \sum_{j=0}^{n-1} B_j.$$

Next, we reduce the difference equation (2) to the simpler equation (3) in the following clever manner. Let

$$z_n = \frac{y_n}{A^n}, \quad n=0, 1, \dots, N.$$

Then, $z_{n+1} = y_{n+1}/A^{n+1}$. But $y_{n+1} = Ay_n + B$. Consequently,

$$z_{n+1} = \frac{y_n}{A^n} + \frac{B}{A^{n+1}} = z_n + \frac{B}{A^{n+1}}.$$

Therefore,

$$\begin{aligned} z_n &= z_0 + \sum_{j=0}^{n-1} \frac{B}{A^{j+1}} = y_0 + \frac{B}{A} \left[\frac{1 - \left(\frac{1}{A}\right)^n}{1 - \frac{1}{A}} \right] \\ &= y_0 + \frac{B}{A-1} \left[1 - \left(\frac{1}{A}\right)^n \right] \end{aligned}$$

and

$$y_n = A^n z_n = A^n y_0 + \frac{B}{A-1} (A^n - 1). \quad (4)$$

Finally, returning to the inequalities (1), we see that

$$E_n \leq \frac{B}{A-1} (A^n - 1), \quad n=1, 2, \dots, N. \quad (5)$$

While collecting material for this book, this author was approached by a colleague with the following problem. He had just received a bill from the bank for the first payment on his wife's student loan. This loan was to be repaid in 10 years in 120 equal monthly installments. According to his rough estimate, the bank was overcharging him by at least 20%. Before confronting the bank's officers, though, he wanted to compute exactly the monthly payments due on this loan.

This problem can be put in the following more general framework. Suppose that P dollars are borrowed from a bank at an annual interest rate of $R\%$. This loan is to be repaid in n years in equal monthly installments of x dollars. Find x .

Our first step in solving this problem is to compute the interest due on the loan. To this end observe that the interest I_1 owed when the first payment is due is $I_1 = (r/12)P$, where $r = R/100$. The principal outstanding during the second month of the loan is $(x - I_1)$ less than the principal outstanding during the first month. Hence, the interest I_2 owed during the second month of the loan is

$$I_2 = I_1 - \frac{r}{12}(x - I_1).$$

Similarly, the interest I_{j+1} owed during the $(j+1)$ st month is

1 First-order differential equations

$$I_{j+1} = I_j - \frac{r}{12}(x - I_j) = \left(1 + \frac{r}{12}\right)I_j - \frac{r}{12}x, \quad (6)$$

where I_j is the interest owed during the j th month.

Equation (6) is a difference equation for the numbers

$$I_1 = \frac{r}{12}P, I_2, \dots, I_{12n}.$$

Its solution (see Exercise 4) is

$$I_j = \frac{r}{12}P \left(1 + \frac{r}{12}\right)^{j-1} + x \left[1 - \left(1 + \frac{r}{12}\right)^{j-1}\right]$$

Hence, the total amount of interest paid on the loan is

$$\begin{aligned} I &= I_1 + I_2 + \dots + I_{12n} = \sum_{j=1}^{12n} I_j \\ &= \frac{r}{12}P \sum_{j=1}^{12n} \left(1 + \frac{r}{12}\right)^{j-1} + 12nx - x \sum_{j=1}^{12n} \left(1 + \frac{r}{12}\right)^{j-1} \end{aligned}$$

Now,

$$\sum_{j=1}^{12n} \left(1 + \frac{r}{12}\right)^{j-1} = \frac{12}{r} \left[\left(1 + \frac{r}{12}\right)^{12n} - 1 \right].$$

Therefore,

$$\begin{aligned} I &= P \left[\left(1 + \frac{r}{12}\right)^{12n} - 1 \right] + 12nx - \frac{12x}{r} \left[\left(1 + \frac{r}{12}\right)^{12n} - 1 \right] \\ &= 12nx - P + P \left(1 + \frac{r}{12}\right)^{12n} - \frac{12x}{r} \left[\left(1 + \frac{r}{12}\right)^{12n} - 1 \right]. \end{aligned}$$

But, $12nx - P$ must equal I , since $12nx$ is the amount of money paid the bank and P was the principal loaned. Consequently,

$$P \left(1 + \frac{r}{12}\right)^{12n} - \frac{12x}{r} \left[\left(1 + \frac{r}{12}\right)^{12n} - 1 \right] = 0$$

and this equation implies that

$$x = \frac{\frac{r}{12}P \left(1 + \frac{r}{12}\right)^{12n}}{\left(1 + \frac{r}{12}\right)^{12n} - 1}. \quad (7)$$

Epilog. Using Equation (7), this author computed x for his wife's and his colleague's wife's student loans. In both cases the bank was right—to the penny.

EXERCISES

1. Solve the difference equation $y_{n+1} = -7y_n + 2, y_0 = 1$.

2. Find y_{37} if $y_{n+1} = 3y_n + 1, y_0 = 0, n = 0, 1, \dots, 36$.

3. Estimate the numbers E_0, E_1, \dots, E_N if $E_0 = 0$ and

(a) $E_{n+1} \leq 3E_n + 1, n = 0, 1, \dots, N-1$;

(b) $E_{n+1} \leq 2E_n + 2, n = 0, 1, \dots, N-1$.

4. (a) Show that the transformation $y_j = I_{j+1}$ transforms the difference equation

$$I_{j+1} = \left(1 + \frac{r}{12}\right)I_j - \frac{r}{12}x, \quad I_1 = \frac{r}{12}P$$

into the difference equation

$$y_{j+1} = \left(1 + \frac{r}{12}\right)y_j - \frac{r}{12}x, \quad y_0 = \frac{r}{12}P.$$

(b) Use Equation (4) to find $y_{j-1} = I_j$.

5. Solve the difference equation $y_{n+1} = a_n y_n + b_n, y_1 = \alpha$. *Hint:* Set $z_1 = y_1$ and $z_n = y_n / a_1 \dots a_{n-1}$ for $n \geq 2$. Observe that

$$\begin{aligned} z_{n+1} &= \frac{y_{n+1}}{a_1 \dots a_n} = \frac{a_n y_n}{a_1 \dots a_n} + \frac{b_n}{a_1 \dots a_n} \\ &= z_n + \frac{b_n}{a_1 \dots a_n}. \end{aligned}$$

Hence, conclude that $z_n = z_1 + \sum_{j=1}^{n-1} b_j / a_1 \dots a_j$.

6. Solve the difference equation $y_{n+1} - ny_n = 1 - n, y_1 = 2$.

7. Find y_{25} if $y_1 = 1$ and $(n+1)y_{n+1} - ny_n = 2^n, n = 1, \dots, 24$.

8. A student borrows P dollars at an annual interest rate of $R\%$. This loan is to be repayed in n years in equal monthly installments of x dollars. Find x if

(a) $P = 4250, R = 3$, and $n = 5$;

(b) $P = 5000, R = 7$, and $n = 10$.

9. A home buyer takes out a \$30,000 mortgage at an annual interest rate of 9%. This loan is to be repayed over 20 years in 240 equal monthly installments of x dollars.

(a) Compute x .

(b) Find x if the annual interest rate is 10%.

10. The quantity supplied of some commodity in a given week is obviously an increasing function of its price the previous week, while the quantity demanded in a given week is a function of its current price. Let S_j and D_j denote, respectively, the quantities supplied and demanded in the j th week, and let P_j denote the price of the commodity in the j th week. We assume that there exist positive constants a, b , and c such that

$$S_j = aP_{j-1} \quad \text{and} \quad D_j = b - cP_j.$$

(a) Show that $P_j = b/(a+c) + (-a/c)^j (P_0 - b/(a+c))$, if supply always equals demand.

1 First-order differential equations

- (b) Show that P_j approaches $b/(a+c)$ as j approaches infinity if $a/c < 1$.
- (c) Show that $P = b/(a+c)$ represents an equilibrium situation. That is to say, if supply always equals demand, and if the price ever reaches the level $b/(a+c)$, then it will always remain at that level.

1.13 Numerical approximations; Euler's method

In Section 1.9 we showed that it is not possible, in general, to solve the initial-value problem

$$\frac{dy}{dt} = f(t, y), \quad y(t_0) = y_0. \quad (1)$$

Therefore, in order that differential equations have any practical value for us, we must devise ways of obtaining accurate approximations of the solution $y(t)$ of (1). In Sections 1.13–1.16 we will derive algorithms, which can be implemented on a digital computer, for obtaining accurate approximations of $y(t)$.

Now, a computer obviously cannot approximate a function on an entire interval $t_0 \leq t \leq t_0 + a$ since this would require an infinite amount of information. At best it can compute approximate values y_1, \dots, y_N of $y(t)$ at a finite number of points t_1, t_2, \dots, t_N . However, this is sufficient for our purpose since we can use the numbers y_1, \dots, y_N to obtain an accurate approximation of $y(t)$ on the entire interval $t_0 \leq t \leq t_0 + a$. To wit, let $\hat{y}(t)$ be the function whose graph on each interval $[t_j, t_{j+1}]$ is the straight line connecting the points (t_j, y_j) and (t_{j+1}, y_{j+1}) (see Figure 1). We can express $\hat{y}(t)$ analytically by the equation

$$\hat{y}(t) = y_j + \frac{1}{h}(t - t_j)(y_{j+1} - y_j), \quad t_j \leq t \leq t_{j+1}.$$

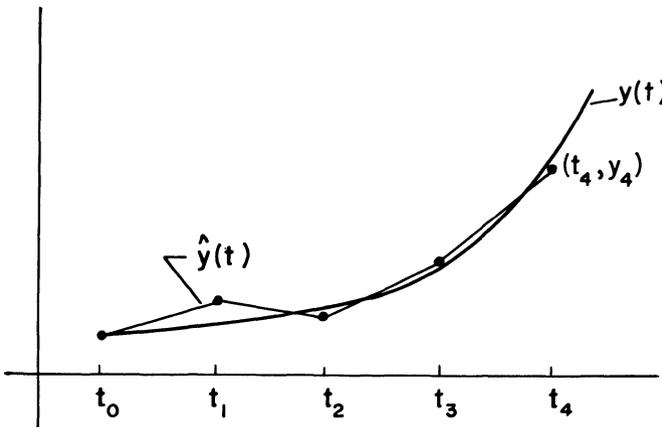


Figure 1. Comparison of $\hat{y}(t)$ and $y(t)$

If $\hat{y}(t)$ is close to $y(t)$ at $t = t_j$; that is, if y_j is close to $y(t_j)$, and if t_{j+1} is close to t_j , then $\hat{y}(t)$ is close to $y(t)$ on the entire interval $t_j \leq t \leq t_{j+1}$. This follows immediately from the continuity of both $y(t)$ and $\hat{y}(t)$. Thus, we need only devise schemes for obtaining accurate approximations of $y(t)$ at a discrete number of points t_1, \dots, t_N in the interval $t_0 \leq t \leq t_0 + a$. For simplicity, we will require that the points t_1, \dots, t_N be equally spaced. This is achieved by choosing a large integer N and setting $t_k = t_0 + k(a/N)$, $k = 1, \dots, N$. Alternately, we may write $t_{k+1} = t_k + h$ where $h = a/N$.

Now, the only thing we know about $y(t)$ is that it satisfies a certain differential equation, and that its value at $t = t_0$ is y_0 . We will use this information to compute an approximate value y_1 of y at $t = t_1 = t_0 + h$. Then, we will use this approximate value y_1 to compute an approximate value y_2 of y at $t = t_2 = t_1 + h$, and so on. In order to accomplish this we must find a theorem which enables us to compute the value of y at $t = t_k + h$ from the knowledge of y at $t = t_k$. This theorem, of course, is Taylor's Theorem, which states that

$$y(t_k + h) = y(t_k) + h \frac{dy(t_k)}{dt} + \frac{h^2}{2!} \frac{d^2y(t_k)}{dt^2} + \dots \quad (2)$$

Thus, if we know the value of y and its derivatives at $t = t_k$, then we can compute the value of y at $t = t_k + h$. Now, $y(t)$ satisfies the initial-value problem (1). Hence, its derivative, when evaluated at $t = t_k$, must equal $f(t_k, y(t_k))$. Moreover, by repeated use of the chain rule of partial differentiation (see Appendix A), we can evaluate

$$\frac{d^2y(t_k)}{dt^2} = \left[\frac{\partial f}{\partial t} + f \frac{\partial f}{\partial y} \right] (t_k, y(t_k))$$

and all other higher-order derivatives of $y(t)$ at $t = t_k$. Hence, we can re-write (2) in the form

$$\begin{aligned} y(t_{k+1}) &= y(t_k) + hf(t_k, y(t_k)) \\ &\quad + \frac{h^2}{2!} \left[\frac{\partial f}{\partial t} + f \frac{\partial f}{\partial y} \right] (t_k, y(t_k)) + \dots \end{aligned} \quad (3)$$

The simplest approximation of $y(t_{k+1})$ is obtained by truncating the Taylor series (3) after the second term. This gives rise to the numerical scheme

$$y_1 = y_0 + hf(t_0, y_0), \quad y_2 = y_1 + hf(t_1, y_1),$$

and, in general,

$$y_{k+1} = y_k + hf(t_k, y_k), \quad y_0 = y(t_0). \quad (4)$$

Notice how we use the initial-value y_0 and the fact that $y(t)$ satisfies the differential equation $dy/dt = f(t, y)$ to compute an approximate value y_1 of $y(t)$ at $t = t_1$. Then, we use this approximate value y_1 to compute an approximate value y_2 of $y(t)$ at $t = t_2$, and so on.

1 First-order differential equations

Equation (4) is known as *Euler's scheme*. It is the simplest numerical scheme for obtaining approximate values y_1, \dots, y_N of the solution $y(t)$ at times t_1, \dots, t_N . Of course, it is also the least accurate scheme, since we have only retained two terms in the Taylor series expansion for $y(t)$. As we shall see shortly, Euler's scheme is not accurate enough to use in many problems. However, it is an excellent introduction to the more complicated schemes that will follow.

Example 1. Let $y(t)$ be the solution of the initial-value problem

$$dy/dt = 1 + (y - t)^2, \quad y(0) = \frac{1}{2}.$$

Use Euler's scheme to compute approximate values y_1, \dots, y_N of $y(t)$ at the points $t_1 = 1/N, t_2 = 2/N, \dots, t_N = 1$.

Solution. Euler's scheme for this problem is

$$y_{k+1} = y_k + h[1 + (y_k - t_k)^2], \quad k = 0, 1, \dots, N-1, \quad h = 1/N$$

with $y_0 = \frac{1}{2}$. Sample Pascal and Fortran programs to compute y_1, \dots, y_N are given below. These programs, as well as all subsequent programs, have variable values for t_0, y_0, a , and N , so that they may also be used to solve the more general initial-value problem $dy/dt = 1 + (y - t)^2, y(t_0) = y_0$ on any desired interval. Moreover, these same programs work even if we change the differential equation; if we change the function $f(t, y)$ then we need only change line 12 in the Pascal program (and line 11 in the C program) and the expressions for $Y(1)$ and $Y(K)$ in Section B of the Fortran program.

Pascal Program

Program Euler (input, output);

var

T, Y: array[0..999] of real;

a, h: real;

k, N: integer;

begin

readln(T[0], Y[0], a, N);

h := a/N;

page;

for k := 0 to N - 1 do

begin

T[k + 1] := T[k] + h;

Y[k + 1] := Y[k] + h * (1 + (Y[k] - T[k]) * (Y[k] - T[k]));

end;

writeln('T':4, 'Y':16);

for k := 0 to N do

writeln(T[k]:10:7, ' ':2, Y[k]:16:9);

end.

Fortran Program

<p>Section A Read in data</p>	}	10		<pre>DIMENSION T(1000), Y(1000) READ (5, 10) T0, Y0, A, N FORMAT (3F20.8, I5) H=A/N</pre>
<p>Section B Do computations</p>	}	20		<pre>T(1)=T0+H Y(1)=Y0+H*(1+(Y0-T0)**2) DO 20 K=2, N T(K)=T(K-1)+H Y(K)=Y(K-1)+H*(1+(Y(K-1) -T(K-1))**2) CONTINUE</pre>
<p>Section C Print out results</p>	}	30		<pre>WRITE (6, 30) T0, Y0, (T(J), Y(J), J=1, N) FORMAT (1H1, 3X, 1HT, 4X, 1HY, / (1H, 1X, F10.7, 2X, F20.9/)) CALL EXIT END</pre>

See also C Program 4 in Appendix C for a sample C program. Table 1 below gives the results of these computations for $a = 1$, $N = 10$, $t_0 = 0$, and $y_0 = \frac{1}{2}$. All of these computations, and all subsequent computations, were carried out on an IBM 360 computer using 16 decimal places accuracy. The results have been rounded to 8 significant decimal places.

Table 1

t	y	t	y
0	0.5	0.6	1.29810115
0.1	0.625	0.7	1.44683567
0.2	0.7525625	0.8	1.60261202
0.3	0.88309503	0.9	1.76703063
0.4	1.01709501	1	1.94220484
0.5	1.15517564		

The exact solution of this initial-value problem (see Exercise 7) is

$$y(t) = t + 1/(2 - t).$$

Thus, the error we make in approximating the value of the solution at $t = 1$ by y_{10} is approximately 0.06, since $y(1) = 2$. If we run this program for $N = 20$ and $N = 40$, we obtain that $y_{20} = 1.96852339$ and $y_{40} = 1.9835109$. Hence, the error we make in approximating $y(1)$ by y_{40} is already less than 0.02.

EXERCISES

Using Euler's method with step size $h=0.1$, determine an approximate value of the solution at $t=1$ for each of the initial-value problems 1–5. Repeat these computations with $h=0.025$ and compare the results with the given value of the solution.

1. $\frac{dy}{dt} = 1 + t - y, \quad y(0)=0; (y(t)=t)$

2. $\frac{dy}{dt} = 2ty, \quad y(0)=2; (y(t)=2e^{t^2})$

3. $\frac{dy}{dt} = 1 + y^2 - t^2, \quad y(0)=0; (y(t)=t)$

4. $\frac{dy}{dt} = te^{-y} + \frac{t}{1+t^2}, \quad y(0)=0; (y(t)=\ln(1+t^2))$

5. $\frac{dy}{dt} = -1 + 2t + \frac{y^2}{(1+t^2)^2}, \quad y(0)=1; (y(t)=1+t^2)$

6. Using Euler's method with $h=\pi/40$, determine an approximate value of the solution of the initial-value problem

$$\frac{dy}{dt} = 2\sec^2 t - (1+y^2), \quad y(0)=0$$

at $t=\pi/4$. Repeat these computations with $h=\pi/160$ and compare the results with the number one which is the value of the solution $y(t)=\tan t$ at $t=\pi/4$.

7. (a) Show that the substitution $y=t+z$ reduces the initial-value problem $y'=1+(y-t)^2, y(0)=0.5$ to the simpler initial-value problem $z'=z^2, z(0)=0.5$.
 (b) Show that $z(t)=1/(2-t)$. Hence, $y(t)=t+1/(2-t)$.

1.13.1 Error analysis for Euler's method

One of the nice features of Euler's method is that it is relatively simple to estimate the error we make in approximating $y(t_k)$ by y_k . Unfortunately, though, we must make the severe restriction that t_1, \dots, t_N do not exceed $t_0 + \alpha$, where α is the number defined in the existence–uniqueness theorem of Section 1.10. More precisely, let a and b be two positive numbers and assume that the functions $f, \partial f/\partial t$, and $\partial f/\partial y$ are defined and continuous in the rectangle $t_0 \leq t \leq t_0 + a, y_0 - b \leq y \leq y_0 + b$. We will denote this rectangle by R . Let M be the maximum value of $|f(t,y)|$ for (t,y) in R , and set $\alpha = \min(a, b/M)$. We will determine the error committed in approximating $y(t_k)$ by y_k , for $t_k \leq t_0 + \alpha$.

To this end observe that the numbers y_0, y_1, \dots, y_N satisfy the difference equation

$$y_{k+1} = y_k + hf(t_k, y_k), \quad k=0, 1, \dots, N-1 \quad (1)$$

while the numbers $y(t_0), y(t_1), \dots, y(t_N)$ satisfy the difference equation

$$y(t_{k+1}) = y(t_k) + hf(t_k, y(t_k)) + \frac{h^2}{2} \left[\frac{\partial f}{\partial t} + f \frac{\partial f}{\partial y} \right] (\xi_k, y(\xi_k)) \quad (2)$$

where ξ_k is some number between t_k and t_{k+1} . Equation (2) follows from the identity

$$\frac{d^2 y}{dt^2} = \frac{\partial f}{\partial t} + f \frac{\partial f}{\partial y}$$

and the fact that

$$y(t+h) = y(t) + h \frac{dy(t)}{dt} + \frac{h^2}{2} \frac{d^2 y(\tau)}{dt^2},$$

for some number τ between t and $t+h$. Subtracting Equation (1) from Equation (2) gives

$$\begin{aligned} y(t_{k+1}) - y_{k+1} &= y(t_k) - y_k + h[f(t_k, y(t_k)) - f(t_k, y_k)] \\ &\quad + \frac{h^2}{2} \left[\frac{\partial f}{\partial t} + f \frac{\partial f}{\partial y} \right] (\xi_k, y(\xi_k)). \end{aligned}$$

Next, observe that

$$f(t_k, y(t_k)) - f(t_k, y_k) = \frac{\partial f(t_k, \eta_k)}{\partial y} [y(t_k) - y_k]$$

where η_k is some number between $y(t_k)$ and y_k . Consequently,

$$\begin{aligned} |y(t_{k+1}) - y_{k+1}| &\leq |y(t_k) - y_k| + h \left| \frac{\partial f(t_k, \eta_k)}{\partial y} \right| |y(t_k) - y_k| \\ &\quad + \frac{h^2}{2} \left| \left[\frac{\partial f}{\partial t} + f \frac{\partial f}{\partial y} \right] (\xi_k, y(\xi_k)) \right|. \end{aligned}$$

In order to proceed further, we must obtain estimates of the quantities $(\partial f(t_k, \eta_k))/\partial y$ and $[(\partial f/\partial t) + f(\partial f/\partial y)](\xi_k, y(\xi_k))$. To this end observe that the points $(\xi_k, y(\xi_k))$ and (t_k, y_k) all lie in the rectangle R . (It was shown in Section 1.10 that the points $(\xi_k, y(\xi_k))$ lie in R . In addition, a simple induction argument (see Exercise 9) shows that the points (t_k, y_k) all lie in R .) Consequently, the points (t_k, η_k) must also lie in R . Let L and D be two positive numbers such that

$$\max_{(t,y) \text{ in } R} \left| \frac{\partial f}{\partial y} \right| \leq L$$

and

$$\max_{(t,y) \text{ in } R} \left| \frac{\partial f}{\partial t} + f \frac{\partial f}{\partial y} \right| \leq D.$$

Such numbers always exist if f , $\partial f/\partial t$, and $\partial f/\partial y$ are continuous in R . Then,

$$|y(t_{k+1}) - y_{k+1}| \leq |y(t_k) - y_k| + hL|y(t_k) - y_k| + \frac{Dh^2}{2}. \quad (3)$$

1 First-order differential equations

Now, set $E_k = |y(t_k) - y_k|$, $k=0, 1, \dots, N$. The number E_k is the error we make at the k th step in approximating $y(t_k)$ by y_k . From (3)

$$E_{k+1} \leq (1 + hL)E_k + \frac{Dh^2}{2}, \quad k=0, 1, \dots, N-1. \quad (4)$$

Moreover, $E_0 = 0$ since $y(t_0) = y_0$. Thus, the numbers E_0, E_1, \dots, E_N satisfy the set of inequalities

$$E_{k+1} \leq AE_k + B, \quad E_0 = 0$$

with $A = 1 + hL$ and $B = Dh^2/2$. Consequently, (see Section 1.12)

$$E_k \leq \frac{B}{A-1}(A^k - 1) = \frac{Dh}{2L} [(1 + hL)^k - 1]. \quad (5)$$

We can also obtain an estimate for E_k that is independent of k . Observe that $1 + hL \leq e^{hL}$. This follows from the fact that

$$\begin{aligned} e^{hL} &= 1 + hL + \frac{(hL)^2}{2!} + \frac{(hL)^3}{3!} + \dots \\ &= (1 + hL) + \text{“something positive”}. \end{aligned}$$

Therefore,

$$E_k \leq \frac{Dh}{2L} [(e^{hL})^k - 1] = \frac{Dh}{2L} [e^{khL} - 1].$$

Finally, since $kh \leq \alpha$, we see that

$$E_k \leq \frac{Dh}{2L} [e^{\alpha} - 1], \quad k = 1, \dots, N. \quad (6)$$

Equation (6) says that the error we make in approximating the solution $y(t)$ at time $t = t_k$ by y_k is at most a fixed constant times h . This suggests, as a rule of thumb, that our error should decrease by approximately $\frac{1}{2}$ if we decrease h by $\frac{1}{2}$. We can verify this directly in Example 1 of the previous section where our error at $t = 1$ for $h = 0.1, 0.05$, and 0.025 is $0.058, 0.032$, and 0.017 respectively.

Example 1. Let $y(t)$ be the solution of the initial-value problem

$$\frac{dy}{dt} = \frac{t^2 + y^2}{2}, \quad y(0) = 0.$$

- (a) Show that $y(t)$ exists at least for $0 \leq t \leq 1$, and that in this interval, $-1 \leq y(t) \leq 1$.
 (b) Let N be a large positive integer. Set up Euler's scheme to find approximate values of y at the points $t_k = k/N$, $k = 0, 1, \dots, N$.
 (c) Determine a step size $h = 1/N$ so that the error we make in approximating $y(t_k)$ by y_k does not exceed 0.0001 .

Solution. (a) Let R be the rectangle $0 \leq t \leq 1$, $-1 \leq y \leq 1$. The maximum value that $(t^2 + y^2)/2$ achieves for (t, y) in R is 1 . Hence, by the exist-

ence-uniqueness theorem of Section 1.10, $y(t)$ exists at least for

$$0 \leq t \leq \alpha = \min\left(1, \frac{1}{1}\right) = 1,$$

and in this interval, $-1 \leq y \leq 1$.

$$(b) \quad y_{k+1} = y_k + h \left(\frac{t_k^2 + y_k^2}{2} \right) = y_k + \frac{1}{2N} \left[\left(\frac{k}{N} \right)^2 + y_k^2 \right]$$

with $y_0 = 0$. The integer k runs from 0 to $N - 1$.

(c) Let $f(t, y) = (t^2 + y^2)/2$, and compute

$$\frac{\partial f}{\partial y} = y \quad \text{and} \quad \frac{\partial f}{\partial t} + f \frac{\partial f}{\partial y} = t + \frac{y}{2}(t^2 + y^2).$$

From (6), $|y(t_k) - y_k| \leq (Dh/2L)(e^L - 1)$ where L and D are two positive numbers such that

$$\max_{(t, y) \in R} |y| \leq L$$

and

$$\max_{(t, y) \in R} \left| t + \frac{y}{2}(t^2 + y^2) \right| \leq D.$$

Now, the maximum values of the functions $|y|$ and $|t + (y/2)(t^2 + y^2)|$ for (t, y) in R are clearly 1 and 2 respectively. Hence,

$$|y(t_k) - y_k| \leq \frac{2h}{2}(e - 1) = h(e - 1).$$

This implies that the step size h should be smaller than $0.0001/(e - 1)$. Equivalently, N should be larger than $(e - 1)10^4 = 17,183$. Thus, we must iterate the equation

$$y_{k+1} = y_k + \frac{1}{2(17,183)} \left[\left(\frac{k}{17,183} \right)^2 + y_k^2 \right]$$

17,183 times to be sure that $y(1)$ is correct to four decimal places.

Example 2. Let $y(t)$ be the solution of the initial-value problem

$$\frac{dy}{dt} = t^2 + e^{-y^2}, \quad y(0) = 1.$$

(a) Show that $y(t)$ exists at least for $0 \leq t \leq 1$, and that in this interval, $-1 \leq y \leq 3$.

(b) Let N be a large positive integer. Set up Euler's scheme to find approximate values of $y(t)$ at the points $t_k = k/N$, $k = 0, 1, \dots, N$.

(c) Determine a step size h so that the error we make in approximating $y(t_k)$ by y_k does not exceed 0.0001.

Solution.

(a) Let R be the rectangle $0 \leq t \leq 1$, $|y - 1| \leq 2$. The maximum value that $t^2 + e^{-y^2}$ achieves for (t, y) in R is 2. Hence, $y(t)$ exists at least for $0 \leq t \leq \min(1, 2/2) = 1$, and in this interval, $-1 \leq y \leq 3$.

1 First-order differential equations

(b) $y_{k+1} = y_k + h(t_k^2 + e^{-y_k^2}) = y_k + (1/N)[(k/N)^2 + e^{-y_k^2}]$ with $y_0 = 1$. The integer k runs from 0 to $N - 1$.

(c) Let $f(t, y) = t^2 + e^{-y^2}$ and compute

$$\frac{\partial f}{\partial y} = -2ye^{-y^2}, \quad \text{and} \quad \frac{\partial f}{\partial t} + f \frac{\partial f}{\partial y} = 2t - 2y(t^2 + e^{-y^2})e^{-y^2}.$$

From (6), $|y(t_k) - y_k| \leq (Dh/2L)(e^L - 1)$ where L and D are two positive numbers such that

$$\max_{(t, y) \text{ in } R} |-2ye^{-y^2}| \leq L$$

and

$$\max_{(t, y) \text{ in } R} |2t - 2y(t^2 + e^{-y^2})e^{-y^2}| \leq D.$$

Now, it is easily seen that the maximum value of $|2ye^{-y^2}|$ for $-1 \leq y \leq 3$ is $\sqrt{2/e}$. Thus, we take $L = \sqrt{2/e}$. Unfortunately, though, it is extremely difficult to compute the maximum value of the function

$$|2t - 2y(t^2 + e^{-y^2})e^{-y^2}|$$

for (t, y) in R . However, we can still find an acceptable value D by observing that for (t, y) in R ,

$$\begin{aligned} \max |2t - 2y(t^2 + e^{-y^2})e^{-y^2}| &\leq \max |2t| + \max |2y(t^2 + e^{-y^2})e^{-y^2}| \\ &\leq \max |2t| + \max |2ye^{-y^2}| \times \max (t^2 + e^{-y^2}) \\ &= 2 + 2\sqrt{2/e} = 2(1 + \sqrt{2/e}). \end{aligned}$$

Hence, we may choose $D = 2(1 + \sqrt{2/e})$. Consequently,

$$|y(t_k) - y_k| \leq \frac{2(1 + \sqrt{2/e})h[e^{\sqrt{2/e}} - 1]}{2\sqrt{2/e}}.$$

This implies that the step size h must be smaller than

$$\frac{\sqrt{2/e}}{1 + \sqrt{2/e}} \times \frac{0.0001}{e^{\sqrt{2/e}} - 1}.$$

Examples 1 and 2 show that Euler's method is not very accurate since approximately 20,000 iterations are required to achieve an accuracy of four decimal places. One obvious disadvantage of a scheme which requires so many iterations is the cost. The going rate for computer usage at present is about \$1200.00 per hour. A second, and much more serious disadvantage, is that y_k may be very far away from $y(t_k)$ if N is exceptionally large. To wit, a digital computer can never perform a computation exactly since it only retains a finite number of decimal places. Consequently, every time we perform an arithmetic operation on the computer, we must introduce a

“round off” error. This error, of course, is small. However, if we perform too many operations then the accumulated round off error may become so large as to make our results meaningless. Exercise 8 gives an illustration of this for Euler's method.

EXERCISES

1. Determine an upper bound on the error we make in using Euler's method with step size h to find an approximate value of the solution of the initial-value problem

$$\frac{dy}{dt} = \frac{t^2 + y^2}{2}, \quad y(0) = 1$$

at any point t in the interval $[0, \frac{2}{5}]$. *Hint:* Let R be the rectangle $0 \leq t \leq 1$, $0 \leq y \leq 2$.

2. Determine an upper bound on the error we make in using Euler's method with step size h to find an approximate value of the solution of the initial-value problem

$$\frac{dy}{dt} = t - y^4, \quad y(0) = 0$$

at any point t in the interval $[0, 1]$. *Hint:* Let R be the rectangle $0 \leq t \leq 1$, $-1 \leq y \leq 1$.

3. Determine an upper bound on the error we make in using Euler's method with step size h to find an approximate value of the solution of the initial-value problem

$$\frac{dy}{dt} = t + e^y, \quad y(0) = 0$$

at any point t in the interval $[0, 1/(e+1)]$. *Hint:* Let R be the rectangle $0 \leq t \leq 1$, $-1 \leq y \leq 1$.

4. Determine a suitable value of h so that the error we make in using Euler's method with step size h to find an approximate value of the solution of the initial-value problem

$$\frac{dy}{dt} = e^t - y^2, \quad y(0) = 0$$

at any point t in the interval $[0, 1/e]$ is at most 0.0001. *Hint:* Let R be the rectangle $0 \leq t \leq 1$, $-1 \leq y \leq 1$.

5. Determine a suitable value of h so that the error we make in using Euler's method with step size h to find an approximate value of the solution of the initial-value problem

$$\frac{dy}{dt} = t^2 + \tan^2 y, \quad y(0) = 0$$

at any point t in the interval $[0, \frac{1}{2}]$ is at most 0.00001. *Hint:* Let R be the rectangle $0 \leq t \leq \frac{1}{2}$, $-\pi/4 \leq y \leq \pi/4$.

1 First-order differential equations

6. Determine a suitable value of h so that the error we make in using Euler's method with step size h to find an approximate value of the solution of the initial-value problem

$$\frac{dy}{dt} = \frac{1}{1+t^2+y^2}, \quad y(0)=0$$

at any point t in the interval $[0, 1]$ is at most 0.0001. *Hint:* Let R be the rectangle $0 \leq t \leq 1, -1 \leq y \leq 1$.

7. Let $y(t)$ be the solution of the initial-value problem

$$\frac{dy}{dt} = f(t, y), \quad y(0) = 0.$$

Suppose that $|f(t, y)| \leq 1$, $|\partial f / \partial y| \leq 1$, and $|(\partial f / \partial t) + f(\partial f / \partial y)| \leq 2$ in the rectangle $0 \leq t \leq 1, -1 \leq y \leq 1$. When the Euler scheme

$$y_{k+1} = y_k + hf(t_k, y_k), \quad h = \frac{1}{N}$$

is used with $N = 10$, the value of y_5 is $-0.15[(\frac{11}{10})^5 - 1]$, and the value of y_6 is $0.12[(\frac{11}{10})^6 - 1]$. Prove that $y(t)$ is zero at least once in the interval $(\frac{1}{2}, \frac{3}{5})$.

8. Let $y(t)$ be the solution of the initial-value problem

$$y' = f(t, y), \quad y(t_0) = y_0.$$

Euler's method for finding approximate values of $y(t)$ is $y_{k+1} = y_k + hf(t_k, y_k)$. However, the quantity $y_k + hf(t_k, y_k)$ is never computed exactly: we always introduce an error ϵ_k with $|\epsilon_k| < \epsilon$. That is to say, the computer computes numbers $\tilde{y}_1, \tilde{y}_2, \dots$, such that

$$\tilde{y}_{k+1} = \tilde{y}_k + hf(t_k, \tilde{y}_k) + \epsilon_k$$

with $\tilde{y}_0 = y_0$. Suppose that $|\partial f / \partial y| \leq L$ and $|(\partial f / \partial t) + f(\partial f / \partial y)| \leq D$ for all t and y .

- (a) Show that

$$E_{k+1} \equiv |y(t_{k+1}) - \tilde{y}_{k+1}| \leq (1 + hL)E_k + \frac{D}{2}h^2 + \epsilon$$

- (b) Conclude from (a) that

$$E_k \leq \left[\frac{Dh}{2} + \frac{\epsilon}{h} \right] \frac{e^{\alpha L} - 1}{L}$$

for $kh \leq \alpha$.

- (c) Choose h so that the error E_k is minimized. Notice that the error E_k may be very large if h is very small.

9. Let y_1, y_2, \dots satisfy the recursion relation

$$y_{k+1} = y_k + hf(t_k, y_k).$$

Let R be the rectangle $t_0 \leq t \leq t_0 + a, y_0 - b \leq y \leq y_0 + b$, and assume that $|f(t, y)| \leq M$ for (t, y) in R . Finally, let $\alpha = \min(a, b/M)$.

- (a) Prove that $|y_j - \tilde{y}_0| \leq jhM$, as long as $jh \leq \alpha$. *Hint:* Use induction.

- (b) Conclude from (a) that the points (t_j, y_j) all lie in R as long as $j \leq \alpha/h$.

1.14 The three term Taylor series method

Euler's method was derived by truncating the Taylor series

$$y(t_{k+1}) = y(t_k) + hf(t_k, y(t_k)) + \frac{h^2}{2} \left[\frac{\partial f}{\partial t} + f \frac{\partial f}{\partial y} \right] (t_k, y(t_k)) + \dots \quad (1)$$

after the second term. The most obvious way of obtaining better numerical schemes is to retain more terms in Equation (1). If we truncate this Taylor series after three terms then we obtain the numerical scheme

$$y_{k+1} = y_k + hf(t_k, y_k) + \frac{h^2}{2} \left[\frac{\partial f}{\partial t} + f \frac{\partial f}{\partial y} \right] (t_k, y_k), \quad k=0, \dots, N-1 \quad (2)$$

with $y_0 = y(t_0)$.

Equation (2) is called the *three term Taylor series method*. It is obviously more accurate than Euler's method. Hence, for fixed h , we would expect that the numbers y_k generated by Equation (2) are better approximations of $y(t_k)$ than the numbers y_k generated by Euler's scheme. This is indeed the case, for it can be shown that $|y(t_k) - y_k|$ is proportional to h^2 whereas the error we make using Euler's method is only proportional to h . The quantity h^2 is much less than h if h is very small. Thus, the three term Taylor series method is a significant improvement over Euler's method.

Example 1. Let $y(t)$ be the solution of the initial-value problem

$$\frac{dy}{dt} = 1 + (y - t)^2, \quad y(0) = \frac{1}{2}.$$

Use the three term Taylor series method to compute approximate values of $y(t)$ at the points $t_k = k/N$, $k = 1, \dots, N$.

Solution. Let $f(t, y) = 1 + (y - t)^2$. Then,

$$\frac{\partial f}{\partial t} + f \frac{\partial f}{\partial y} = -2(y - t) + 2(y - t)[1 + (y - t)^2] = 2(y - t)^3.$$

Hence, the three term Taylor series scheme is

$$y_{k+1} = y_k + h[1 + (y_k - t_k)^2] + h^2(y_k - t_k)^3$$

with $h = 1/N$ and $y_0 = \frac{1}{2}$. The integer k runs from 0 to $N - 1$. Sample Pascal and Fortran programs to compute y_1, \dots, y_N are given below. Again, these programs have variable values for t_0 , y_0 , a , and N .

1 First-order differential equations

Pascal Program

```
Program Taylor (input, output);
var
  T, Y: array[0..999] of real;
  a, h, Temp: real;
  k, N: integer;
begin
  readln(T[0], Y[0], a, N);
  h := a/N;
  page;
  for k := 0 to N - 1 do
    begin
      Temp := Y[k] - T[k];
      T[k + 1] := T[k] + h;
      Y[k + 1] := Y[k] + h * (1 + Temp * Temp)
        + h * h * Temp * Temp * Temp;
    end;
  writeln('T':4, 'Y':16);
  for k := 0 to N do
    writeln(T[k]:10:7, ' ' :2, Y[k]:16:9);
  end.
```

Fortran Program

Replace Section B of the Fortran program in Section 1.13 by the following:

```
20 | T(1)=T0+H
   | D2Y=H*(Y0-T0)**3
   | Y(1)=Y0+H*(D2Y+1+(Y0-T0)**2)
   | D0 20 K=2,N
   | T(K)=T(K-1)+H
   | D2Y=H*(Y(K-1)-T(K-1))**3
   | Y(K)=Y(K-1)+H*(D2Y+1+(Y(K-1)-T(K-1))**2)
   | CONTINUE
```

See also C Program 5 in Appendix C for a sample C program.

Table 1 below shows the results of these computations for $a = 1$, $N = 10$, $t_0 = 0$, and $y_0 = 0.5$.

Now Euler's method with $N = 10$ predicted a value of 1.9422 for $y(1)$. Notice how much closer the number 1.9957 is to the correct value 2. If we run this program for $N = 20$ and $N = 40$, we obtain that $y_{20} = 1.99884247$ and

Table 1

t	y	t	y
0	0.5	0.6	1.31331931
0.1	0.62625	0.7	1.4678313
0.2	0.7554013	0.8	1.63131465
0.3	0.88796161	0.9	1.80616814
0.4	1.02456407	1	1.99572313
0.5	1.1660084		

$y_{40} = 1.99969915$. These numbers are also much more accurate than the values 1.96852339 and 1.9835109 predicted by Euler's method.

EXERCISES

Using the three term Taylor series method with $h=0.1$, determine an approximate value of the solution at $t=1$ for each of the initial-value problems 1–5. Repeat these computations with $h=0.025$ and compare the results with the given value of the solution.

- $dy/dt = 1 + t - y$, $y(0) = 0$; ($y(t) = t$)
- $dy/dt = 2ty$, $y(0) = 2$; ($y(t) = 2e^{t^2}$)
- $dy/dt = 1 + y^2 - t^2$, $y(0) = 0$; ($y(t) = t$)
- $dy/dt = te^{-y} + t/(1+t^2)$, $y(0) = 0$; ($y(t) = \ln(1+t^2)$)
- $dy/dt = -1 + 2t + y^2/(1+t^2)^2$, $y(0) = 1$; ($y(t) = 1 + t^2$)
- Using the three term Taylor series method with $h = \pi/40$, determine an approximate value of the solution of the initial-value problem

$$\frac{dy}{dt} = 2\sec^2 t - (1 + y^2), \quad y(0) = 0$$

at $t = \pi/4$. Repeat these computations with $h = \pi/160$ and compare the results with the number one which is the value of the solution $y(t) = \tan t$ at $t = \pi/4$.

1.15 An improved Euler method

The three term Taylor series method is a significant improvement over Euler's method. However, it has the serious disadvantage of requiring us to compute partial derivatives of $f(t, y)$, and this can be quite difficult if the function $f(t, y)$ is fairly complicated. For this reason we would like to derive numerical schemes which do not require us to compute partial derivatives of $f(t, y)$. One approach to this problem is to integrate both sides of the differential equation $y' = f(t, y)$ between t_k and $t_k + h$ to obtain that

$$y(t_{k+1}) = y(t_k) + \int_{t_k}^{t_k+h} f(t, y(t)) dt. \quad (1)$$

1 First-order differential equations

This reduces the problem of finding an approximate value of $y(t_{k+1})$ to the much simpler problem of approximating the area under the curve $f(t, y(t))$ between t_k and $t_k + h$. A crude approximation of this area is $hf(t_k, y(t_k))$, which is the area of the rectangle R in Figure 1a. This gives rise to the numerical scheme

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

which, of course, is Euler's method.

A much better approximation of this area is

$$\frac{h}{2} [f(t_k, y(t_k)) + f(t_{k+1}, y(t_{k+1}))]$$

which is the area of the trapezoid T in Figure 1b. This gives rise to the numerical scheme

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

However, we cannot use this scheme to determine y_{k+1} from y_k since y_{k+1} also appears on the right-hand side of (2). A very clever way of overcoming this difficulty is to replace y_{k+1} in the right-hand side of (2) by the value $y_k + hf(t_k, y_k)$ predicted for it by Euler's method. This gives rise to the numerical scheme

$$y_{k+1} = y_k + \frac{h}{2} [f(t_k, y_k) + f(t_k + h, y_k + hf(t_k, y_k))], \quad y_0 = y(t_0). \quad (3)$$

Equation (3) is known as the *improved Euler method*. It can be shown that $|y(t_k) - y_k|$ is at most a fixed constant times h^2 . Hence, the improved Euler method gives us the same accuracy as the three term Taylor series method without requiring us to compute partial derivatives.

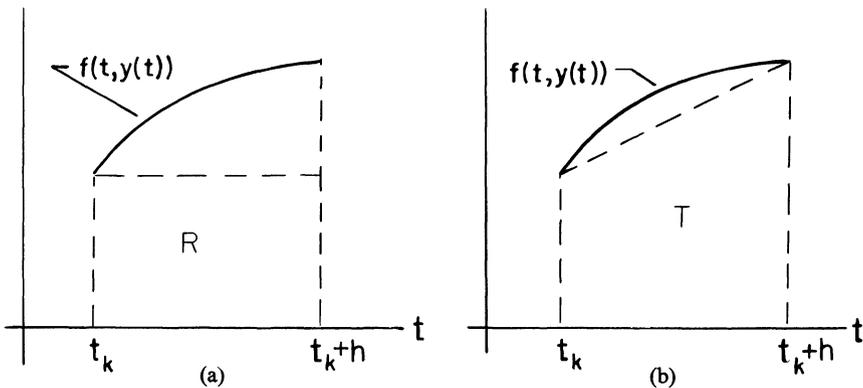


Figure 1

Example 1. Let $y(t)$ be the solution of the initial-value problem

$$\frac{dy}{dt} = 1 + (y - t)^2, \quad y(0) = \frac{1}{2}.$$

Use the improved Euler method to compute approximate values of $y(t)$ at the points $t_k = k/N$, $k = 1, \dots, N$.

Solution. The improved Euler scheme for this problem is

$$y_{k+1} = y_k + \frac{h}{2} \left\{ 1 + (y_k - t_k)^2 + 1 + \left[y_k + h(1 + (y_k - t_k)^2) - t_{k+1} \right]^2 \right\}$$

with $h = 1/N$ and $y_0 = 0.5$. The integer k runs from 0 to $N - 1$. Sample Pascal and Fortran programs to compute y_1, \dots, y_N are given below. Again, these programs have variable values for t_0, y_0, a , and N .

Pascal Program

Program Improved (input, output);

var

T, Y: array[0..999] of real;

a, h, R: real;

k, N: integer;

begin

readln(T[0], Y[0], a, N);

h := a/N;

page;

for k:=0 to N-1 do

begin

R := 1 + (Y[k] - T[k]) * (Y[k] - T[k]);

T[k+1] := T[k] + h;

Y[k+1] := Y[k] + (h/2) * (R + 1

+ (Y[k] + h * R - T[k+1]) * (Y[k] + h * R - T[k+1]));

end;

writeln('T':4, 'Y':16);

for k := 0 to N do

writeln(T[k]:10:7, ' ':2, Y[k]:16:9);

end.

Fortran Program

Replace Section B of the Fortran program in Example 1 of Section 1.13 by the following:

```

20 | T(1)=T0+H
   | R=1+(Y0-T0)**2
   | Y(1)=Y0+(H/2)*(R+1+(Y0+(H*R)-T(1))**2)
   | DO 20 K=2,N
   | T(K)=T(K-1)+H
   | R=1+(Y(K-1)-T(K-1))**2
   | Y(K)=Y(K-1)+(H/2)*(R+1+(Y(K-1)+(H*R)-T(K))**2)
   | CONTINUE

```

See also C Program 6 in Appendix C for a sample C program.

Table 1 below shows the results of these computations for $a=1$, $N=10$, $t_0=0$, and $y_0=0.5$. If we run this program for $N=20$ and $N=40$ we obtain that $y_{20}=1.99939944$ and $y_{40}=1.99984675$. Hence the values y_{10} , y_{20} , and y_{40} computed by the improved Euler method are even closer to the correct value 2 than the corresponding values 1.99572313, 1.99884246, and 1.99969915 computed by the three term Taylor series method.

Table 1

t	y	t	y
0	0.5	0.6	1.31377361
0.1	0.62628125	0.7	1.46848715
0.2	0.75547445	0.8	1.63225727
0.3	0.88809117	0.9	1.80752701
0.4	1.02477002	1	1.99770114
0.5	1.16631867		

EXERCISES

Using the improved Euler method with $h=0.1$, determine an approximate value of the solution at $t=1$ for each of the initial-value problems 1–5. Repeat these computations with $h=0.025$ and compare the results with the given value of the solution.

- $dy/dt=1+t-y$, $y(0)=0$; ($y(t)=t$)
- $dy/dt=2ty$, $y(0)=2$; ($y(t)=2e^{t^2}$)
- $dy/dt=1+y^2-t^2$, $y(0)=0$; ($y(t)=t$)
- $dy/dt=te^{-y}+t/(1+t^2)$, $y(0)=0$; ($y(t)=\ln(1+t^2)$)
- $dy/dt=-1+2t+y^2/(1+t^2)^2$, $y(0)=1$; ($y(t)=1+t^2$)
- Using the improved Euler method with $h=\pi/40$, determine an approximate value of the solution of the initial-value problem

$$\frac{dy}{dt}=2\sec^2 t-(1+y^2), \quad y(0)=0$$

at $t=\pi/4$. Repeat these computations with $h=\pi/160$ and compare the results with the number one which is the value of the solution $y(t)=\tan t$ at $t=\pi/4$.

1.16 The Runge–Kutta method

We now present, without proof, a very powerful scheme which was developed around 1900 by the mathematicians Runge and Kutta. Because of its simplicity and great accuracy, the Runge–Kutta method is still one of the most widely used numerical schemes for solving differential equations. It is defined by the equation

$$y_{k+1} = y_k + \frac{h}{6} [L_{k,1} + 2L_{k,2} + 2L_{k,3} + L_{k,4}], \quad k = 0, 1, \dots, N-1$$

where $y_0 = y(t_0)$ and

$$\begin{aligned} L_{k,1} &= f(t_k, y_k), & L_{k,2} &= f\left(t_k + \frac{1}{2}h, y_k + \frac{1}{2}hL_{k,1}\right) \\ L_{k,3} &= f\left(t_k + \frac{1}{2}h, y_k + \frac{1}{2}hL_{k,2}\right), & L_{k,4} &= f(t_k + h, y_k + hL_{k,3}). \end{aligned}$$

This formula involves a weighted average of values of $f(t, y)$ taken at different points. Hence the sum $\frac{1}{6}[L_{k,1} + 2L_{k,2} + 2L_{k,3} + L_{k,4}]$ can be interpreted as an average slope. It can be shown that the error $|y(t_k) - y_k|$ is at most a fixed constant times h^4 . Thus, the Runge–Kutta method is much more accurate than Euler's method, the three term Taylor series method and the improved Euler method.

Example 1. Let $y(t)$ be the solution of the initial-value problem

$$\frac{dy}{dt} = 1 + (y - t)^2, \quad y(0) = \frac{1}{2}.$$

Use the Runge–Kutta method to find approximate values y_1, \dots, y_N of y at the points $t_k = k/N$, $k = 1, \dots, N$.

Solution. Sample Pascal and Fortran programs to compute y_1, \dots, y_N by the Runge–Kutta method are given below. These programs differ from our previous programs in that they do not compute y_1 separately. Rather, they compute y_1 in the same “loop” as they compute y_2, \dots, y_N . This is accomplished by relabeling the numbers t_0 and y_0 as t_1 and y_1 respectively.

Pascal Program

Program Runge_Kutta (input, output);

var

T, Y: array[0..999] of real;
a, h, LK1, LK2, LK3, LK4: real;
k, N: integer;

begin

readln(T[0], Y[0], a, N);

h := a/N;

page;

for k := 0 to N - 1 do

begin

T[k + 1] := T[k] + h;

LK1 := 1 + (Y[k] - T[k]) * (Y[k] - T[k]);

LK2 := 1 + ((Y[k] + (h/2) * LK1) - (T[k] + h/2))
* ((Y[k] + (h/2) * LK1) - (T[k] + h/2));

1 First-order differential equations

```

LK3 := 1 + ((Y[k] + (h/2) * LK2) - (T[k] + h/2))
      * ((Y[k] + (h/2) * LK2) - (T[k] + h/2));
LK4 := 1 + ((Y[k] + h * LK3) - (T[k] + h))
      * ((Y[k] + h * LK3) - (T[k] + h));
Y[k + 1] := Y[k] + (h/6) * (LK1 + LK4 + 2 * (LK2 + LK3));
end;
writeln('T':4, 'Y':16);
for k := 0 to N do
  writeln(T[k]:10:7, ' ':2, Y[k]:16:9);
end.

```

Fortran Program

```

10 | DIMENSION T(1000), Y(1000)
    | READ (5, 10) T(1), Y(1), A, N
    | FORMAT (3F20.8, I5)
    | H = A/N
    | D020 K = 1, N
    | T(K + 1) = T(K) + H
    | REAL LK1, LK2, LK3, LK4
    | LK1 = 1 + (Y(K) - T(K)) * * 2
    | LK2 = 1 + ((Y(K) + (H/2) * LK1) - (T(K) + H/2)) * * 2
    | LK3 = 1 + ((Y(K) + (H/2) * LK2) - (T(K) + H/2)) * * 2
    | LK4 = 1 + ((Y(K) + H * LK3) - (T(K) + H)) * * 2
    | Y(K + 1) = Y(K) + (H/6) * (LK1 + LK4 + 2 * (LK2 + LK3))
20 | CONTINUE
    | NA = N + 1
    | WRITE (6, 30) (T(J), Y(J), J = 1, NA)
30 | FORMAT (1H1, 3X, 1HT, 4X, 1HY, / (1H, 1X, F10.7, 2X, F20.9 /))
    | CALL EXIT
    | END

```

See also C Program 7 in Appendix C for a sample C program.

Table 1 below shows the results of these computations for $a = 1$, $N = 10$, $t_0 = 0$, and $y_0 = 0.5$.

Table 1

t	y	t	y
0	0.5	0.6	1.31428555
0.1	0.62631578	0.7	1.4692305
0.2	0.75555536	0.8	1.6333329
0.3	0.88823526	0.9	1.8090902
0.4	1.02499993	1	1.9999988
0.5	1.16666656		

Notice how much closer the number $y_{10} = 1.9999988$ computed by the Runge–Kutta method is to the correct value 2 than the numbers $y_{10} = 1.94220484$, $y_{10} = 1.99572312$, and $y_{10} = 1.99770114$ computed by the Euler, three term Taylor series and improved Euler methods, respectively. If we run this program for $N=20$ and $N=40$, we obtain that $y_{20} = 1.99999992$ and $y_{40} = 2$. Thus, our approximation of $y(1)$ is already correct to eight decimal places when $h=0.025$. Equivalently, we need only choose $N \geq 40$ to achieve eight decimal places accuracy.

To put the accuracy of the various schemes into proper perspective, let us say that we have three different schemes for numerically solving the initial-value problem $dy/dt = f(t, y)$, $y(0) = 0$ on the interval $0 \leq t \leq 1$, and that the error we make in using these schemes is $3h$, $11h^2$, and $42h^4$, respectively. If our problem is such that we require eight decimal places accuracy, then the step sizes h_1 , h_2 , and h_3 of these three schemes must satisfy the inequalities $3h_1 \leq 10^{-8}$, $11h_2^2 \leq 10^{-8}$, and $42h_3^4 \leq 10^{-8}$. Hence, the number of iterations N_1 , N_2 , and N_3 of these three schemes must satisfy the inequalities

$$N_1 \geq 3 \times 10^8 = 300,000,000, \quad N_2 \geq \sqrt{11} \times 10^4 \approx 34,000$$

and

$$N_3 \geq (42)^{1/4} \times 10^2 \approx 260.$$

This is a striking example of the difference between the Runge–Kutta method and the Euler, improved Euler and three term Taylor series methods.

Remark. It should be noted that we perform four functional evaluations at each step in the Runge–Kutta method, whereas we only perform one functional evaluation at each step in Euler’s method. Nevertheless, the Runge–Kutta method still beats the heck out of Euler’s method, the three term Taylor series method, and the improved Euler method.

EXERCISES

Using the Runge–Kutta method with $h = 0.1$, determine an approximate value of the solution at $t = 1$ for each of the initial-value problems 1–5. Repeat these computations with $h = 0.025$ and compare the results with the given value of the solution.

1. $dy/dt = 1 + t - y$, $y(0) = 0$; $(y(t) = t)$
2. $dy/dt = 2ty$, $y(0) = 2$; $(y(t) = 2e^{t^2})$
3. $dy/dt = 1 + y^2 - t^2$, $y(0) = 0$; $(y(t) = t)$
4. $dy/dt = te^{-y} + t/(1 + t^2)$, $y(0) = 0$; $(y(t) = \ln(1 + t^2))$
5. $dy/dt = -1 + 2t + y^2/((1 + t^2)^2)$, $y(0) = 1$; $(y(t) = 1 + t^2)$

1 First-order differential equations

6. Using the Runge–Kutta method with $h = \pi/40$, determine an approximate value of the solution of the initial-value problem

$$\frac{dy}{dt} = 2 \sec^2 t - (1 + y^2), \quad y(0) = 0$$

at $t = \pi/4$. Repeat these computations with $h = \pi/160$ and compare the results with the number one which is the value of the solution $y(t) = \tan t$ at $t = \pi/4$.

1.17 What to do in practice

In this section we discuss some of the practical problems which arise when we attempt to solve differential equations on the computer. First, and foremost, is the problem of estimating the error that we make. It is not too difficult to show that the error we make using Euler's method, the three term Taylor series method, the improved Euler method and the Runge–Kutta method with step size h is at most $c_1 h$, $c_2 h^2$, $c_3 h^2$, and $c_4 h^4$, respectively. With one exception, though, it is practically impossible to find the constants c_1 , c_2 , c_3 , and c_4 . The one exception is Euler's method where we can explicitly estimate (see Section 1.13.1) the error we make in approximating $y(t_k)$ by y_k . However, this estimate is not very useful, since it is only valid for t_k sufficiently close to t_0 , and we are usually interested in the values of y at times t much larger than t_0 . Thus, we usually do not know, a priori, how small to choose the step size h so as to achieve a desired accuracy. We only know that the approximate values y_k that we compute get closer and closer to $y(t_k)$ as h gets smaller and smaller.

One way of resolving this difficulty is as follows. Using one of the schemes presented in the previous section, we choose a step size h and compute numbers y_1, \dots, y_N . We then repeat the computations with a step size $h/2$ and compare the results. If the changes are greater than we are willing to accept, then it is necessary to use a smaller step size. We keep repeating this process until we achieve a desired accuracy. For example, suppose that we require the solution of the initial-value problem $y' = f(t, y)$, $y(0) = y_0$ at $t = 1$ to four decimal places accuracy. We choose a step size $h = 1/100$, say, and compute y_1, \dots, y_{100} . We then repeat these computations with $h = 1/200$ and obtain new approximations z_1, \dots, z_{200} . If y_{100} and z_{200} agree in their first four decimal places then we take z_{200} as our approximation of $y(1)$.* If y_{100} and z_{200} do not agree in their first four decimal places, then we repeat our computations with step size $h = 1/400$.

Example 1. Find the solution of the initial-value problem

*This does not guarantee that z_{200} agrees with $y(1)$ to four decimal places. As an added precaution, we might halve the step size again. If the first four decimal places still remain unchanged, then we can be reasonably certain that z_{200} agrees with $y(1)$ to four decimal places.

$$\frac{dy}{dt} = y(1 + e^{-y}) + e^t, \quad y(0) = 0$$

at $t = 1$ to four decimal places accuracy.

Solution. We illustrate how to try and solve this problem using Euler's method, the three term Taylor series method, the improved Euler method, and the Runge–Kutta method.

(i) *Euler's method:*

Pascal Program

Program Euler (input, output);

var

T, Y: array[0..999] of real;

a, h: real;

k, N: integer;

begin

readln(T[0], Y[0], a, N);

h := a/N;

page;

for k := 0 to N - 1 do

begin

T[k + 1] := T[k] + h;

Y[k + 1] := Y[k] + h * (Y[k] * (1 + exp(-Y[k])) + exp(T[k]));

end;

writeln('N':4, 'h':10, 'Y[N]':20);

writeln(N:4, ' ':2, h:10:7, ' ':2, Y[N]:18:10);

end.

Fortran Program

Section A Read in data	}	10		DIMENSION T(1000), Y(1000) READ (5, 10) T(1), Y(1), A, N, FORMAT (3F20.8, I5) H = A/N
Section B Do computations	}	20	1	DO 20 K = 1, N T(K + 1) = T(K) + H Y(K + 1) = Y(K) + H * (Y(K) * (1 + EXP(-Y(K))) + EXP(T(K))) CONTINUE
Section C Print out results	}	30		WRITE (6, 30) N, H, Y(N + 1) FORMAT (1H, 1X, I5, 2X, F10.7, 4X, F20.9) CALL EXIT END

1 First-order differential equations

See also C Program 8 in Appendix C for a sample C program. We set $A = 1$, $T[0] = 0$, $Y[0] = 0$ ($T(1) = Y(1) = 0$ in the Fortran program) and ran these programs for $N = 10, 20, 40, 80, 160, 320$, and 640 . The results of these computations are given in Table 1. Notice that even with a step size h

Table 1

N	h	y_N
10	0.1	2.76183168
20	0.05	2.93832741
40	0.025	3.03202759
80	0.0125	3.08034440
160	0.00625	3.10488352
320	0.003125	3.11725009
640	0.0015625	3.12345786

as small as $1/640$, we can only guarantee an accuracy of one decimal place. This points out the limitation of Euler's method. Since N is so large already, it is wiser to use a more accurate scheme than to keep choosing smaller and smaller step sizes h for Euler's method.

(ii) *The three term Taylor series method*

Pascal Program

```
Program Taylor (input, output);
var
  T, Y: array[0..999] of real;
  a, h, DY1, DY2: real;
  k, N: integer;
begin
  readln(T[0], Y[0], a, N);
  h := a/N;
  page;
  for k:=0 to N-1 do
    begin
      T[k+1] := T[k] + h;
      DY1 := 1 + (1 - Y[k]) * exp(-Y[k]);
      DY2 := Y[k] * (1 + exp(-Y[k])) + exp(T[k]);
      Y[k+1] := Y[k] + h * DY2 + (h * h/2) * (exp(T[k]) + DY1 * DY2);
    end;
  writeln('N':4, 'h':10, 'Y[N]':20);
  writeln(N:4, ' ':2, h:10:7, ' ':2, Y[N]:18:10);
end.
```

Fortran Program

Replace Section B of the previous Fortran program by

```

20 | D020 K=1,N
    | T(K+1)=T(K)+H
    | DY1=1+(1-Y(K))*EXP(-Y(K))
    | DY2=Y(K)*(1+EXP(-Y(K)))+EXP(T(K))
    | Y(K+1)=Y(K)+H*DY2+(H*H/2)*(EXP(T(K))+DY1*DY2)
    | CONTINUE

```

See also C Program 9 in Appendix C for a sample C program.

We set $A=1$, $T[0]=0$, and $Y[0]=0$ ($T(1)=0$, and $Y(1)=0$ in the Fortran program) and ran these programs for $N=10, 20, 40, 60, 80, 160$, and 320 . The results of these computations are given in Table 2. Observe that y_{160}

Table 2

N	h	y_N
10	0.1	3.11727674
20	0.05	3.12645293
40	0.025	3.12885845
80	0.0125	3.12947408
160	0.00625	3.12962979
320	0.003125	3.12966689

and y_{320} agree in their first four decimal places. Hence the approximation $y(1)=3.12966689$ is correct to four decimal places.

(iii) *The improved Euler method*

Pascal Program

Program Improved (input, output);

var

T, Y: array[0..999] of real;
a, h, R1, R2: real;
k, N: integer;

begin

 readln(T[0], Y[0], a, N);
 h:=a/N;
 page;
 for k:=0 to N-1 do
 begin
 T[k+1]:=T[k]+h;

1 First-order differential equations

```

R1 := Y[k] * (1 + exp(-Y[k])) + exp(T[k]);
R2 := (Y[k] + h * R1) * (1 + exp(-(Y[k] + h * R1))) + exp(T[k + 1]);
Y[k + 1] := Y[k] + (h/2) * (R1 + R2);
end;
writeln('N':4, 'h':10, 'Y[N]':20);
writeln(N:4, ' ':2, h:10:7, ' ':2, Y[N]:18:10);
end.

```

Fortran Program

Replace Section B of the first Fortran program in this section by

```

20 | D020 K = 1, N
    | T(K + 1) = T(K) + H
    | R1 = Y(K) * (1 + EXP(-Y(K))) + EXP(T(K))
    | R2 = (Y(K) + H * R1) * (1 + EXP(-(Y(K) + H * R1))) + EXP(T(K + 1))
    | Y(K + 1) = Y(K) + (H/2) * (R1 + R2)
    | CONTINUE

```

See also C Program 10 in Appendix C for a sample C program.

We set $A = 1$, $T[0] = 0$ and $Y[0] = 0$ ($T(1) = 0$ and $Y(1) = 0$ in the Fortran program) and ran these programs for $N = 10, 20, 40, 80, 160$, and 320 . The results of these computations are given in Table 3. Observe that y_{160} and y_{320}

Table 3

N	h	y_N
10	0.1	3.11450908
20	0.05	3.12560685
40	0.025	3.1286243
80	0.0125	3.12941247
160	0.00625	3.12961399
320	0.003125	3.12964943

agree in their first four decimal places. Hence the approximation $y(1) = 3.12964943$ is correct to four decimal places.

(iv) *The method of Runge–Kutta*

Pascal Program

Program Runge_Kutta (input, output);

```

var
  T, Y: array[0..999] of real;
  a, h, LK1, LK2, LK3, LK4: real;
  k, N: integer;

```

```

begin
  readln(T[0], Y[0], a, N);
  h := a/N;
  page;
  for k := 0 to N - 1 do
    begin
      T[k + 1] := T[k] + h;
      LK1 := (Y[k] * (1 + exp(-Y[k])) + exp(T[k]));
      LK2 := (Y[k] + (h/2) * LK1) * (1 + exp(-(Y[k] + (h/2) * LK1)))
        + exp(T[k] + (h/2));
      LK3 := (Y[k] + (h/2) * LK2) * (1 + exp(-(Y[k] + (h/2) * LK2)))
        + exp(T[k] + (h/2));
      LK4 := (Y[k] + h * LK3) * (1 + exp(-(Y[k] + h * LK3)))
        + exp(T[k + 1]);
      Y[k + 1] := Y[k] + (h/6) * (LK1 + 2 * LK2 + 2 * LK3 + LK4);
    end;
  writeln('N':4, 'h':10, 'Y[N]':20);
  writeln(N:4, ' ':2, h:10:7, ' ':2, Y[N]:18:10);
0 0 1 10

```

Fortran Program

Replace Section B of the first Fortran program in this section by

```

D020 K = 1, N
      T(K + 1) = T(K) + H
      LK1 = Y(K) * (1 + EXP(-Y(K)) + EXP(T(K)))
      LK2 = (Y(K) + (H/2) * LK1) * (1 + EXP(-(Y(K) + (H/2) * LK1)))
1     + EXP(T(K) + (H/2))
      LK3 = (Y(K) + (H/2) * LK2) * (1 + EXP(-(Y(K) + (H/2) * LK2)))
1     + EXP(T(K) + (H/2))
      LK4 = (Y(K) + H * LK3) * (1 + EXP(-(Y(K) + H * LK3)))
1     + EXP(T(K + 1))
20    Y(K + 1) = Y(K) + (H/6) * (LK1 + 2 * LK2 + 2 * LK3 + LK4)
      CONTINUE

```

See also C Program 11 in Appendix C for a sample C program.

We set $A = 1$, $T[0] = 0$ and $Y[0] = 0$ ($T(1) = 0$ and $Y(1) = 0$ in the Fortran program) and ran these programs for $N = 10, 20, 40, 80, 160$, and 320 . The results of these computations are given in Table 4. Notice that our approximation of $y(1)$ is already correct to four decimal places with $h = 0.1$, and that it is already correct to eight decimal places with $h = 0.00625$ ($N = 160$). This example again illustrates the power of the Runge–Kutta method.

We conclude this section with two examples which point out some additional difficulties which may arise when we solve initial-value problems on a digital computer.

1 First-order differential equations

Table 4

N	h	y_N
10	0.1	3.1296517
20	0.05	3.12967998
40	0.025	3.1296819
80	0.0125	3.12968203
160	0.00625	3.12968204
320	0.003125	3.12968204

Example 2. Use the Runge–Kutta method to find approximate values of the solution of the initial-value problem

$$\frac{dy}{dt} = t^2 + y^2, \quad y(0) = 1$$

at the points $t_k = k/N$, $k = 1, \dots, N$.

Solution.

Pascal Program

Program Runge_Kutta (input, output);

var

T, Y: array[0..999] of real;
a, h, LK1, LK2, LK3, LK4: real;
k, N: integer;

begin

 readln(T[0], Y[0], a, N);

 h := a/N;

 page;

 for k := 0 to N - 1 do

 begin

 T[k + 1] := T[k] + h;

 LK1 := T[k] * T[k] + Y[k] * Y[k];

 LK2 := (T[k] + h/2) * (T[k] + h/2)

 + (Y[k] + (h/2) * LK1) * (Y[k] + (h/2) * LK1);

 LK3 := (T[k] + h/2) * (T[k] + h/2)

 + (Y[k] + (h/2) * LK2) * (Y[k] + (h/2) * LK2);

 LK4 := (T[k] + h) * (T[k] + h)

 + (Y[k] + h * LK3) * (Y[k] + h * LK3);

 Y[k + 1] := Y[k] + (h/6) * (LK1 + 2 * LK2 + 2 * LK3 + LK4);

 end;

 writeln('T':4, 'Y':15);

 for k := 0 to N do

 writeln(T[k]:10:7, ' ':2, Y[k]:16:9);

end.

Fortran Program

Replace Sections B and C of the first Fortran program in this section by

Section B Do computa- tions	}	20	D020 K=1,N T(K+1)=T(K)+H LK1=T(K)**2+Y(K)**2 LK2=(T(K)+(H/2))**2+(Y(K)+(H/2)*LK1)**2 LK3=(T(K)+(H/2))**2+(Y(K)+(H/2)*LK2)**2 LK4=(T(K)+H)**2+(Y(K)+H*LK3)**2 Y(K+1)=Y(K)+(H/6)*(LK1+2*LK2+2*LK3+LK4) CONTINUE
Section C Print out results	}	30	NA=N+1 WRITE (6,30) (T(J),Y(J),J=1,NA) FORMAT (1H1,3X,1HT,4X,1HY/(1H,1X,F9.7, 2X,F20.9/)) CALL EXIT END

See also C Program 12 in Appendix C for a sample C program.

We attempted to run these programs with $A=1$, $T[0]=0$, $Y[0]=1$ ($T(1)=0$, and $Y(1)=1$ in the Fortran program) and $N=10$, but we received an error message that the numbers being computed exceeded the domain of the computer. That is to say, they were larger than 10^{38} . This indicates that the solution $y(t)$ goes to infinity somewhere in the interval $[0, 1]$. We can prove this analytically, and even obtain an estimate of where $y(t)$ goes to infinity, by the following clever argument. Observe that for $0 \leq t \leq 1$, $y(t)$ is never less than the solution $\phi_1(t) = 1/(1-t)$ of the initial-value problem

$$\frac{dy}{dt} = y^2, \quad y(0) = 1.$$

In addition, $y(t)$ never exceeds the solution $\phi_2(t) = \tan(t + \pi/4)$ of the initial-value problem $dy/dt = 1 + y^2$, $y(0) = 1$. Hence, for $0 \leq t \leq 1$,

$$\frac{1}{1-t} \leq y(t) \leq \tan(t + \pi/4).$$

This situation is described graphically in Figure 1. Since $\phi_1(t)$ and $\phi_2(t)$ become infinite at $t=1$ and $t=\pi/4$ respectively, we conclude that $y(t)$ becomes infinite somewhere between $\pi/4$ and 1.

The solutions of most initial-value problems which arise in physical and biological applications exist for all future time. Thus, we need not be overly concerned with the problem of solutions going to infinity in finite time or the problem of solutions becoming exceedingly large. On the other hand, though, there are several instances in economics where this problem is of paramount importance. In these instances, we are often interested in

1 First-order differential equations

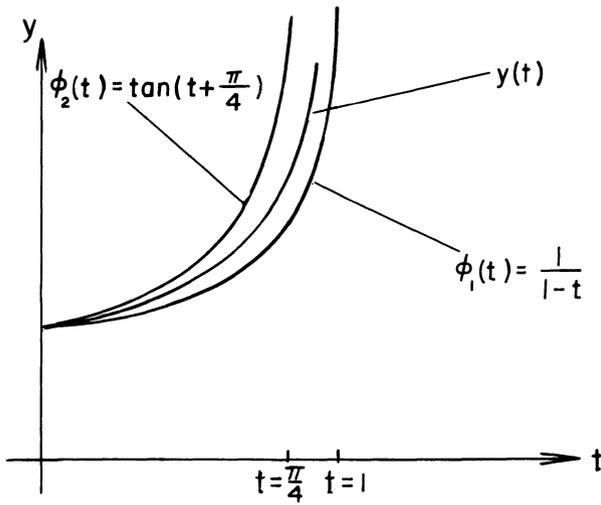


Figure 1

determining whether certain differential equations can accurately model a given economic phenomenon. It is often possible to eliminate several of these equations by showing that they allow solutions which are unrealistically large.

Example 3. Use Euler's method to determine approximate values of the solution of the initial-value problem

$$\frac{dy}{dt} = y|y|^{-3/4} + t \sin \frac{\pi}{t}, \quad y(0) = 0 \quad (1)$$

at the points $1/N, 2/N, \dots, 2$.

Solution. The programming for this problem is simplified immensely by observing that

$$y|y|^{-3/4} = (\text{sgn } y)|y|^{1/4}, \quad \text{where } \text{sgn } y = \begin{cases} 1, & y > 0 \\ 0, & y = 0 \\ -1, & y < 0 \end{cases}$$

Pascal Program

Program Euler (input, output);

const

PI = 3.141592654;

var

T, Y: array[0..999] of real;

h: real;

k, N: integer;

```

begin
  readln(N);
  page;
  h := 2/N;
  T[1] := h;
  Y[1] := 0;
  for k := 1 to N - 1 do
    begin
      T[k + 1] := T[k] + h;
      if Y[k] = 0 then Y[k + 1] := h * T[k] * sin(PI/T[k]) else
        Y[k + 1] := Y[k] + h * (Y[k] * exp((-3/4) * ln(abs(Y[k])))
          + T[k] * sin(PI/T[k]));
    end;
  writeln('T':4, 'Y':15);
  for k := 1 to N do
    writeln(T[k]:10:7, ' ':2, Y[k]:16:9);
end.

```

Fortran Program

```

10 | DIMENSION T(1000), Y(1000)
    | READ (5, 10) N
    | FORMAT (I5)
    | H = 2 / N
    | T(1) = H
    | Y(1) = 0
    | DO 20 K = 2, N
    |   T(K) = T(K - 1) + H
    |   Y(K) = Y(K - 1) + H * (SIGN(Y(K - 1)) * ABS(Y(K - 1)) ** 0.25
    |     + T(K - 1) * SIN(3.141592654 / T(K - 1)))
20 | CONTINUE
    | WRITE(6, 30) 0, 0, (T(J), Y(J), J = 1, N)
30 | FORMAT (1H1, 3X, 1HT, 4X, 1HY / (1H, 1X, F10.7, 2X, F20.9 /))
    | CALL EXIT
    | END

```

See also C Program 13 in Appendix C for a sample C program.

When we set $N = 25$ we obtained the value 2.4844172 for $y(2)$, but when we set $N = 27$, we obtained the value -0.50244575 for $y(2)$. Moreover, all the y_k were positive for $N = 25$ and negative for $N = 27$. We repeated these computations with $N = 89$ and $N = 91$ and obtained the values 2.64286349 and -0.6318074 respectively. In addition, all the y_k were again positive for $N = 89$ and negative for $N = 91$. Indeed, it is possible, but rather difficult, to prove that all the y_k will be positive if $N = 1, 5, 9, 13, 17, \dots$ and negative

1 First-order differential equations

if $N = 3, 7, 11, 15, \dots$. This suggests that the solution of the initial-value problem (1) is not unique. We cannot prove this analytically, since we cannot solve the differential equation explicitly. It should be noted though, that the existence–uniqueness theorem of Section 1.10 does not apply here, since the partial derivative with respect to y of the function $|y|^{-3/4}y + t \sin \pi/t$ does not exist at $y = 0$.

Most of the initial-value problems that arise in applications have unique solutions. Thus, we need not be overly concerned with the problem of non-uniqueness of solutions. However, we should always bear in mind that initial-value problems which do not obey the hypotheses of the existence–uniqueness theorem of Section 1.10 might possess more than one solution, for the consequences of picking the wrong solution in these rare instances can often be catastrophic.

EXERCISES

In each of Problems 1–5, find the solution of the given initial-value problem at $t = 1$ to four decimal places accuracy.

1. $\frac{dy}{dt} = y + e^{-y} + 2t, \quad y(0) = 0$

2. $\frac{dy}{dt} = 1 - t + y^2, \quad y(0) = 0$

3. $\frac{dy}{dt} = \frac{t^2 + y^2}{1 + t + y^2}, \quad y(0) = 0$

4. $\frac{dy}{dt} = e^t y^2 - 2y, \quad y(0) = 1$

5. $\frac{dy}{dt} = ty^3 - y, \quad y(0) = 1$