

Variations on a Theme of Euler*

Robert M. Corless[†]
Julia E. Jankowski[†]

Abstract. This paper is an exploration of numerical methods for solving initial-value problems for ordinary differential equations. We look in detail at one “simple” problem, namely, the leaky bucket, and use it to explore the notions of explicit marching methods vs. implicit marching methods, interpolation, backward error, and stiffness. While the leaky bucket example is very well known, and these topics are studied in a great many textbooks, we will here emphasize backward error in a way that might be new and that we hope will be useful for your students. Indeed, the paper is intended to be a resource for students themselves. We will also use two techniques not normally seen in a first course, a new one, namely, “optimal backward error,” and an old one, namely, “the method of modified equations.”

Key words. Torricelli’s law, numerical methods, modified equations, optimal backward error

AMS subject classification. 65L20

DOI. 10.1137/15M1032351

I. Torricelli’s Law and the Leaky Bucket. Quite a few textbooks use Torricelli’s law, often derived using Bernoulli’s law for steady flow, to write down a differential equation modeling the height of fluid in a bucket with a hole in the bottom. For example, the beautiful book [13]—really one of our favorite textbooks of all time—has an extensive treatment in section 4.2, “Uniqueness: The Leaking Bucket Versus Radioactive Decay.” Its treatment even includes backward error, which they call “slope error.” We (partially) disagree with one thing in the book, however: on p. 172 we find the statement (emphasized in the original) “*numerical methods are inherently untrustworthy, and anything that you might guess using them, especially if it concerns long term behaviour, must be checked by other methods.*” We believe that statement to be only partly true, and this is a good thing too because for many (most?) problems in science nowadays there is *no other way* of obtaining insight (or control or prediction) than by using numerical methods. The purpose of the treatment here is to help learn when to trust numerical methods, and when not to.

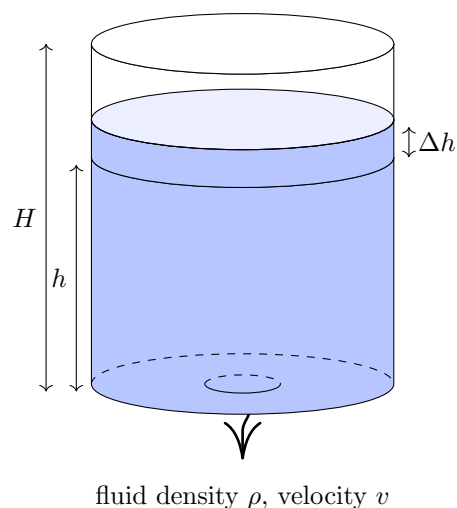
Following [13] we use an energy balance to derive an equation for the height of fluid in a leaky bucket.

The bucket is assumed to be cylindrical with cross-sectional area A , with perfectly vertical walls and a perfectly flat bottom, while the hole in the bottom has an area a . Thus, the volume of fluid lost out of the hole in a time interval $\Delta\tau$ is $a \cdot v(\tau)\Delta\tau$ (if $\Delta\tau$ is small enough that the velocity $v(\tau)$ can be considered constant). At the same

*Received by the editors July 24, 2015; accepted for publication (in revised form) March 22, 2016; published electronically November 3, 2016. This work was supported by NSERC and the Faculty of Science, UWO.

<http://www.siam.org/journals/sirev/58-4/M103235.html>

[†]Department of Applied Mathematics, University of Western Ontario, London, ON, Canada (rcorless@uwo.ca, jjankow4@uwo.ca).

**Fig. 1.1**

time, the volume lost is $A \cdot h'(\tau) \Delta\tau$, where $h'(\tau) = dh/d\tau$; by conservation these are equal, so

$$(1.1) \quad av(\tau) = A \frac{dh}{d\tau}.$$

The change in potential energy in time $\Delta\tau$ is $(\rho A \Delta h)gh$ because the mass is the product of ρ —the density of the fluid—with the volume $A\Delta h$, and the change in energy occurs over the height h because the leak is at the bottom. If there is no energy dissipation (we will revisit this later), then this must equal the kinetic energy of the fluid flowing out, $\frac{1}{2}[\rho A \Delta h]v^2$, plus the increase in kinetic energy of the (slowly) downward flowing fluid in the bucket itself, which like the dissipation we ignore. Thus $\rho A \Delta h gh = \frac{1}{2} \rho A \Delta h v^2$ or

$$(1.2) \quad 2gh = v^2 = \left(\frac{A}{a} \frac{dh}{d\tau} \right)^2.$$

Thus at last we have

$$(1.3) \quad \frac{dh}{d\tau} = -\frac{a}{A} \sqrt{2gh},$$

and we take the negative sign because the height is decreasing over time. Let us now nondimensionalize. Put $y = h/H$, where H is the height of the bucket itself. Then $y = 1$ corresponds to a bucket brim full of liquid, while $y = 0$ means the bucket is empty. Dividing (1.3) by H we have

$$(1.4) \quad \frac{dy}{d\tau} = -\frac{a}{A} \sqrt{\frac{2g}{H}} \cdot \sqrt{y}.$$

Now we rescale time. Put

$$(1.5) \quad t = \frac{a}{A} \sqrt{\frac{2g}{H}} \tau$$

and so, by the chain rule, $\frac{dy}{d\tau} = \frac{dt}{d\tau} \frac{dy}{dt}$ and the scale factors cancel, leaving

$$(1.6) \quad \frac{dy}{dt} = -\sqrt{y}.$$

The initial condition is thus $y(0) = y_0$, the initial fraction full, where $0 \leq y_0 \leq 1$. We will almost always take $y_0 = 1$.

Hubbard and West [13] use this example to great effect for studying uniqueness of solutions. Given an empty bucket, there's no way to tell when it was full—if it ever was. This fact is well displayed in the model, as we will see. The reference solution¹ of this separable equation is easy: $y(t) = (\sqrt{y_0} - t/2)^2$.

This solution is valid only as long as $y > 0$. Investigating when $y = 0$ separately we have $\frac{dy}{dt} = 0$ thereafter, so

$$(1.7) \quad y(t) = \begin{cases} (\sqrt{y_0} - t/2)^2, & 0 \leq t < 2\sqrt{y_0}, \\ 0, & 2\sqrt{y_0} \leq t, \end{cases}$$

and by inspection this is continuous, even differentiable, at $T = 2\sqrt{y_0}$. In dimensional terms, the discharge time when $y_0 = 1$ is predicted to be

$$(1.8) \quad \tau_T = \frac{A}{a} \sqrt{\frac{H}{2g}} \cdot 2 = \frac{A}{a} \sqrt{\frac{2H}{g}}.$$

Notice that $y(t) \geq 0$ always, and of course this is necessary for the $\sqrt{y(t)}$ in (1.6) to be real.

All this is fine, and educational. We don't even need numerical methods for this problem, but, for other equations, we do, and we claim it's worthwhile to have a look at what numerical methods do on this nonlinear problem. Before we do, though, notice one more fact about this equation: the differential equation is *not* Lipschitz continuous at $y = 0$, because $\frac{\partial}{\partial y}(-\sqrt{y}) = -1/(2\sqrt{y})$ is infinite there; hence there is no finite L such that $|\sqrt{y_1} - \sqrt{y_2}| \leq L \cdot |y_1 - y_2|$. This means that the classical existence and uniqueness theorems fail to apply; see, e.g., Theorems 5 and 6, p. 28 of [1]. Indeed, this is the whole point of section 4.2 of [13]. Moreover, the hypotheses for some theorems (see, e.g., Theorem 3.4 in [12] or Theorem 203A in [3]) about the convergence of numerical methods therefore *also* must fail to hold. So let's see what happens, in detail.

2. Forward Euler for the Leaky Bucket. Euler's method for $y'(t) = f(t, y)$ marches forward with the iteration $y_{n+1} = y_n + \Delta t_n f(t_n, y_n)$, where $t_{n+1} = t_n + \Delta t_n$. This is intended to give the *skeleton*² $\{(t_n, y_n)\}_{n=0}^N$ of an approximate solution on some time interval $t_0 \leq t \leq t_N$. The sequence of time steps Δt_n is chosen either for convenience (all Δt_n = some fixed Δt), for reliability (take each Δt_n small enough to ensure accuracy), or for efficiency (take each Δt_n as large as possible given constraints on accuracy). For now, take the convenient route and fix $\Delta t_n = \Delta t$. Our iteration, for (1.6), becomes

$$(2.1) \quad y_{n+1} = y_n - \Delta t \sqrt{y_n}.$$

¹We use the phrase "reference solution" instead of the more common "exact solution," because starting in section 4 we will see that everything is an "exact solution" of some relevant equation.

²In older books, numerical methods *only* produce approximate solutions at discrete points: $\{(t_n, y_n)_{n \geq 0}\}$. We now call this just the bare bones of the solution, hence "skeleton." Numerical solutions of ODEs nowadays include interpolants.

Taking $\Delta t = 1$ to start with, we find that a full bucket $y_0 = 1$ empties in one time step: $y_i = 1 - 1 \cdot \sqrt{1} = 0$. It remains empty thereafter. Possibly this time step is too large. Taking $\Delta t = 1/2$ instead, $y_1 = 1 - \frac{1}{2}\sqrt{1} = \frac{1}{2}$ and $y_2 = \frac{1}{2} - \frac{1}{2} \cdot \sqrt{\frac{1}{2}} \doteq 0.146446$, $y_3 \doteq -0.04489$ is negative, and y_4 is complex! This doesn't seem right.

Taking a more artful time step, we choose $\Delta t = \frac{2}{1+\sqrt{5}} \doteq 0.618\dots$. We will see the reason for this choice in a moment. Then $y_0 = 1$ gives

$$(2.2) \quad y_1 = 1 - \frac{2}{1 + \sqrt{5} \cdot 1} = \frac{-1 + \sqrt{5}}{1 + \sqrt{5}} = \frac{4}{(1 + \sqrt{5})^2}$$

and

$$(2.3) \quad y_2 = \frac{4}{(1 + \sqrt{5})^2} - \frac{2}{(1 + \sqrt{5})} \cdot \frac{2}{(1 + \sqrt{5})} = 0.$$

Thus, with this Δt the bucket empties in two steps, at a time $\Delta t + \Delta t = 4/1 + \sqrt{5} = 1.236$, larger than with $\Delta t = 1$.

Another artful choice is $\Delta t = \frac{2}{1+\sqrt{7+2\sqrt{5}}} \doteq 0.45588\dots$. You, the reader, can verify that $y_0 = 1$ leads in *three* steps to $y_3 = 0$ and an empty bucket, in a time $3\Delta t \doteq 1.26764$, still larger.

It turns out that taking Δt smaller doesn't generally help, except temporarily, unless we're artful and choose just those steps such that $y_N = \Delta t^2$ exactly for some lucky N , when $y_{N+1} = 0$ exactly and the bucket empties. For all other Δt , which means almost all Δt , the iteration terminates when some $y_N < (\Delta t)^2$ and then $y_{N+1} = y_N - \Delta t\sqrt{y_N} < 0$. Once any $y_{N+1} < 0$, the next iterate must be complex.

Something important has happened. The forward, explicit Euler scheme seems to have generically failed for this problem. Not, it's true, until the size of the solution is $\mathcal{O}(\Delta t^2)$, and a reasonable human, computing the solution by hand, would stop at this point and declare the bucket empty, especially since taking Δt smaller and doing it again would get us to an even more nearly empty bucket. But dumb computers, on the other hand... We emphasize that it is quite unusual for the explicit Euler method to fail to be able to take a step.³ Even for problems with singular solutions, its tendency is to keep going!

We will revisit this, but first note that we can remove one parameter from the iteration by putting $y_n = f_n(\Delta t)^2$, whence (2.1) becomes

$$(2.4) \quad f_{n+1} = f_n - \sqrt{f_n}.$$

Of course now $f_0 = \frac{y_0}{(\Delta t)^2}$ can be quite large, but this further "nondimensionalization" will be helpful. We now see the source of the "artful" choices for Δt : We look at the preimages of 0 by inverting $f_{n+1} = f_n - \sqrt{f_n}$,

$$(2.5) \quad f_n = \left(\frac{1 + \sqrt{1 + 4f_{n+1}}}{2} \right)^2,$$

giving us $1, \left(\frac{1+\sqrt{5}}{2}\right)^2$ (whence $\Delta t = \frac{2}{1+\sqrt{5}}$), and $\left(\frac{1+\sqrt{1+(1+\sqrt{5})^2}}{2}\right)^2 = \left(\frac{1+\sqrt{7+2\sqrt{5}}}{2}\right)^2$,

³Apparently, Euler knew this could happen.

and so on. For $y_0 = f_0/(\Delta t)^2 = 1$, we see that taking $\Delta t = 1/\sqrt{f_{-N}}$ for any preimage f_{-N} of 0, the bucket empties in a finite number of steps. But for almost all Δt , and certainly if there are rounding errors, we will miss these “sweet spots” and wind up with complex y . Let’s try another method.

Remark 2.1. One reviewer pointed out that we have some freedom in coding the differential equation. We don’t have to limit ourselves to just the equation as it was derived, and we can try to head off any difficulties that we might foresee. Since \sqrt{y} occurs in the equation, we might anticipate trouble if y ever becomes less than 0. It shouldn’t, we know, and we shouldn’t have to worry, but prudent programmers put in safeguards. One way to do this (suggested by the reviewer) is to code the equation as

$$(2.6) \quad \frac{dy}{dt} = \begin{cases} -\sqrt{y} & \text{if } y > 0, \\ 0 & \text{if } y \leq 0. \end{cases}$$

This way at least if y overshoots the equilibrium when the bucket empties only by a small amount, then thereafter the solution will stay constant at that small negative value, which a reasonable human might accept with a wink. However, this did not come from the modeling, but rather from some geometric analysis of the model. Indeed, some studies of Torricelli’s law work from the equation $(\dot{y})^2 + y = 0$ and give a decidedly different analysis [14].

Even if an implicit method is used it needs to avoid negative y in the solution to the implicit equation. If a nonlinear root finder is used in the process of taking a step, it is quite possible that negative y might be sampled in the course of solving the implicit expression. The reformulation discussed would avoid any possible problems with a nonlinear solver here without changing the behavior of backward Euler.

Further, nowadays most software actually used for solving ODEs can handle events such as y reaching zero. It is important for students to know that this sort of issue arises in many practical applications and has been successfully recognized and dealt with by software developers. A better “precautionary coding” then might be to use this feature; that is, instead of integrating over a fixed time interval, one integrates until the event $y = 0$ is detected. But we now return to the more primitive situation and consider another method.

3. Euler Backward for the Leaky Bucket. Implicit Euler, also known as backward Euler or Euler backward, is of course $y_{n+1} = y_n + \Delta t f(t_{n+1}, y_{n+1})$. Here we have

$$(3.1) \quad y_{n+1} = y_n - \Delta t \sqrt{y_{n+1}}.$$

Normally, i.e., for general nonlinear f , this equation must be solved numerically at each step, making the process expensive and worth the bother only when there is something to gain: typically an implicit method is only better than an explicit method when the problem is “stiff.” Indeed, that’s the pragmatic *definition* of a stiff problem: a stiff problem is one for which an implicit method is more effective. For a more mathematical discussion of stiffness, see [16].

Here, though, we can usefully solve the nonlinear equation once and for all and get an explicit iteration again. Put $p = \sqrt{y_{n+1}}$. Notice that $p > 0$. Then $p^2 + \Delta t p - y_n = 0$

is just a quadratic equation and the positive root can be written as

$$(3.2) \quad p = \frac{-\Delta t + \sqrt{(\Delta t^2) + 4y_n}}{2}$$

$$(3.3) \quad = \frac{-\Delta t + \sqrt{(\Delta t^2) + 4y_n}}{2} \cdot \frac{\Delta t + \sqrt{(\Delta t^2) + 4y_n}}{\Delta t + \sqrt{(\Delta t^2) + 4y_n}}$$

$$(3.4) \quad = \frac{-(\Delta t^2) + (\Delta t^2) + 4y_n}{2(\Delta t + \sqrt{(\Delta t^2) + 4y_n})}$$

$$(3.5) \quad = \frac{2y_n}{\Delta t + \sqrt{(\Delta t^2) + 4y_n}},$$

using a standard trick to avoid catastrophic cancellation when y_n is very small [4], so the formula

$$(3.6) \quad y_{n+1} = p^2 = \frac{4y_n^2}{(\Delta t + \sqrt{(\Delta t^2) + 4y_n})^2}$$

lets us iterate implicit Euler in an explicit manner. Of course, this is unusual and only possible because the example is so “simple.” Again we can scale by (Δt^2) to remove a parameter: put $y_n = b_n(\Delta t^2)$ so

$$(3.7) \quad b_{n+1} = \frac{4}{(1 + \sqrt{1 + 4b_n})^2} b_n^2.$$

Taking a few possible initial values, say, $b_0 = 10$ or $b_0 = 1$, gives quite different behavior to forward Euler. Taking $b_0 = 10$ is typical: after only a few iterations,⁴ five in this case, $b_5 \doteq 0.9193$ is less than 1, and thereafter convergence to 0 is quadratically fast—reminiscent of Newton’s method—and indeed $b_{15} \doteq 10^{-302}$ and b_{16} underflows to 0. This behavior is universal as long as $b_0 < \mathcal{O}(1/\varepsilon_m)$, where ε_m is the machine epsilon. No matter what other initial condition $b_0 > 0$ is given, eventually some $b_N < 1$ and thereafter the bucket empties very rapidly. In perfect arbitrary precision arithmetic it would never *actually* empty, but working with values of $y_n = b_n \Delta t^2$ smaller than the reciprocal of Avogadro’s number⁵ seems like overkill for a continuum fluid model.⁶

It seems that Euler backward is better for the leaky bucket. This indicates that the differential equation $\dot{y} = -\sqrt{y}$ is “stiff” near $y = 0$, and indeed it is. But there’s more to say.

4. Optimal Backward Error. Since the 1980s, at least, people have known that one way to assess the error of a numerical method is to interpolate the skeleton (somehow) and compute the residual (also called the defect, the deviation, or the slope error). Specifically and without loss of generality, suppose that the interpolant on $t_n \leq t \leq t_{n+1}$ is called $z_n(t)$; the residual is then

$$(4.1) \quad r_n(t) := z'_n(t) - f(t, z_n(t))$$

⁴If $b_n > 1$, $4b_n/(1 + \sqrt{1 + 4b_n})^2 < 1$ and $b_{n+1} < b_n$; if $b_n \leq 1$, $b_{n+1} < b_n^2$.

⁵Avogadro’s number, 6.0221413×10^{23} , counts the number of molecules in a mole. A 10-liter bucket of water would contain about 55.5 moles, when full.

⁶This discussion echoes a point made by L. N. Trefethen in [17]: the continuum model breaks down before floating point does!

on that same interval. Then the piecewise continuous function $z(t)$ defined by

$$z(t) = \begin{cases} z_0(t), & t_0 \leq t \leq t_1, \\ z_1(t), & t_1 \leq t \leq t_2, \\ \vdots & \\ z_{N-1}(t), & t_{N-1} \leq t \leq t_N, \end{cases}$$

satisfies the differential equation

$$(4.2) \quad z'(t) = f(t, z(t)) + r(t),$$

where the piecewise function $r(t)$ is defined similarly as

$$r(t) = \begin{cases} r_0(t), & t_0 \leq t < t_1, \\ r_1(t), & t_1 \leq t < t_2, \\ \vdots & \\ r_{N-1}(t), & t_{N-1} \leq t \leq t_N, \end{cases}$$

where now $r(t)$ can be discontinuous at the skeleton. This remark bears emphasis. The computed solution $z(t)$ is the exact solution of a *different* differential equation. If $r(t)$ is very small, but the equation $\dot{y} = f(y)$ is somehow perfect and must not be changed, then this observation is only indirectly helpful: one must appeal, e.g., to the Gröbner–Alekseev nonlinear variation of constants formula to say

$$(4.3) \quad y(t) - z(t) = \int_0^t G(t, \tau) r(\tau) d\tau$$

and then bound G in order to bound the forward error [12]. Here, $G(t, \tau)$ is a function that depends on the differential equation and the reference solution y and can in principle be computed a priori, and is in fact a kind of nonlinear Green's function. See [12] for details.

But if, as is generally true in science, the differential equation is only one model among many, then we may already be finished.⁷ Imagine, for instance, the bucket being shaken gently with time: the dynamic effects on the fluid would be quite complicated, but one would be tempted to model them with a small time-dependent forcing term. And anyone who's ever emptied a bottle knows that the flow regime at the end is quite different to what's gone before. Constant velocity v is not very realistic when the bucket is nearly empty! One has to worry about the impact of such forcing and the details of the fluid, e.g., its viscosity—but then, one has to do that anyway in a responsible solution. “Everybody knows this,” but not every textbook or course or software package emphasizes it.

For numerically generated $r(t)$ it's well known that the value of $r(t)$ depends on the quality of the interpolant. One needs high-order accuracy throughout the time step to ensure the residual is small everywhere in the interval. For general solvers, considerable attention is paid to this issue [8]. Indeed, one might wonder whether some other interpolant could produce a smaller $r(t)$, as do Hubbard and West: “It is

⁷Well, we still have to study conditioning or sensitivity, perhaps by using Gröbner–Alekseev. But now the point is robustness, not forward error.

unclear from this derivation whether a wigglier curve joining curve points might even give a higher order dependence on $r(t)$; however, looking at the actual evidence in section 3.3 shows that it won't" [13, p. 183].

We're going to answer these doubts analytically. We will think about an *optimal* interpolant for this problem and compute the smallest possible backward error explicitly. The scope and utility of this technique for general systems $f \in \mathbb{C}^n$ is explored in [6]. Its use for the Dahlquist test problem—which generates some surprises—is explored in [7]. For the present example we do not need a general method of optimization and we can just use the triangle inequality. It's better here to think about *relative* backward error, as well. That is, instead of

$$(4.4) \quad \frac{dz}{dt} = -\sqrt{z(t)} + r(t)$$

we use $\delta(t) = -r(t)/\sqrt{z(t)}$, so $\delta(t)$ is the residual $r(t)$ divided by what the derivative should be, to get

$$(4.5) \quad \begin{aligned} \frac{dz}{dt} &= -\sqrt{z(t)} + r(t) \\ &= -\sqrt{z(t)} - \sqrt{z(t)}\delta(t) \\ &= -\sqrt{z(t)}(1 + \delta(t)), \end{aligned}$$

so $r(t) = -\sqrt{z(t)}\delta(t)$ and

$$(4.6) \quad -\frac{1}{\sqrt{z}} \frac{dz}{dt} - 1 = \delta(t).$$

Integrating both sides over the time step,

$$(4.7) \quad \int_{s=t_n}^{t_{n+1}} -\frac{1}{\sqrt{z(s)}} \dot{z}(s) \, ds - \int_{s=t_n}^{t_{n+1}} 1 \, ds = \int_{s=t_n}^{t_{n+1}} \delta(s) \, ds,$$

so

$$(4.8) \quad 2\sqrt{y_n} - 2\sqrt{y_{n+1}} - \Delta t = \int_{s=t_n}^{t_{n+1}} \delta(s) \, ds$$

because $z(t_n) = y_n$, $z(t_{n+1}) = y_{n+1}$, and $t_{n+1} = t_n + \Delta t$. This constrains how small $\delta(t)$ can be: in other words, it has to be large enough to push the solution to y_{n+1} . One can show from this that if y_{n+1} was generated from y_n by a p th-order method, then $\|\delta\| = \mathcal{O}(\Delta t^p)$. For Euler's method, $p = 1$.

Explicitly, for Euler's method we have $y_{n+1} = y_n - \Delta t\sqrt{y_n}$ and so

$$(4.9) \quad \begin{aligned} 2\sqrt{y_n} - 2\sqrt{y_{n+1}} - \Delta t &= 2\sqrt{y_n} - 2\sqrt{y_n - \Delta t\sqrt{y_n}} - \Delta t \\ &= 2\sqrt{y_n} - 2\sqrt{y_n} \left(1 - \frac{\Delta t}{\sqrt{y_n}}\right)^{\frac{1}{2}} - \Delta t \\ &= 2\sqrt{y_n} - 2\sqrt{y_n} \left(1 - \frac{\Delta t}{2\sqrt{y_n}} - \frac{1}{8} \frac{\Delta t^2}{y_n} + \mathcal{O}(\Delta t^3)\right) - \Delta t \\ &= 2\sqrt{y_n} - 2\sqrt{y_n} + \Delta t - \Delta t + \frac{1}{4\sqrt{y_n}} \Delta t^2 + \mathcal{O}(\Delta t^3) \\ &= \frac{1}{4\sqrt{y_n}} \Delta t^2 + \mathcal{O}(\Delta t^3). \end{aligned}$$

Since

$$(4.10) \quad \int_0^{\Delta t} \delta(s) \, ds = \bar{\delta} \Delta t,$$

where $\bar{\delta}$ is the average value of $\delta(s)$, we have

$$(4.11) \quad \bar{\delta} \Delta t = \frac{1}{4\sqrt{y_n}} \Delta t^2 + \mathcal{O}(\Delta t^3)$$

or

$$(4.12) \quad \bar{\delta} = \frac{1}{4\sqrt{y_n}} \Delta t + \mathcal{O}(\Delta t^2).$$

That is, on average, $\delta(s)$ must be $\mathcal{O}(\Delta t^p)$, with $p = 1$. This means that the maximum δ must be at least $\mathcal{O}(\Delta t)$; we will now see that we can find an interpolant such that δ is exactly this size.

From (4.8), the triangle inequality gives

$$(4.13) \quad |2\sqrt{y_n} - 2\sqrt{y_{n+1}} - \Delta t| = \left| \int_{t_n}^{t_{n+1}} \delta(s) \, ds \right| \leq \Delta t \cdot \|\delta(t)\|_\infty$$

and equality is attained when $\delta(t)$ is constant, namely, $\delta(t) \equiv \delta_n$, where $2\sqrt{y_n} - 2\sqrt{y_{n+1}} - \Delta t = \Delta t \delta_n$ or

$$(4.14) \quad \delta_n = \frac{2}{\Delta t} (\sqrt{y_n} - \sqrt{y_{n+1}}) - 1.$$

For clarity and emphasis we're going to repeat ourselves. The triangle inequality proves that

$$(4.15) \quad \|\delta\|_\infty \geq \left| \frac{2}{\Delta t} (\sqrt{y_n} - \sqrt{y_{n+1}}) - 1 \right|,$$

giving a lower bound on how big $\delta(t)$ must be to have $dz/dt = -\sqrt{z}(1 + \delta(t))$, starting with $z(t_n) = y_n$ and making it all the way to $z(t_{n+1}) = y_{n+1}$. Remember, y_{n+1} is determined by the numerical method.

Moreover, this lower bound is achieved if $\delta(t) = 2(\sqrt{y_n} - \sqrt{y_{n+1}})/\Delta t - 1$ is constant; that is, this $\delta(t)$ achieves transit from y_n to y_{n+1} with the least possible relative deviation from the original model. Thinking about backward error in this way is likely new to you. We hope to convince you it will be valuable for your students.

Remark 4.1. This solution works *only* for this problem as *long* as $z(t)$ is never zero in the interval; in particular, it fails when $y_{n+1} < 0$. For that, we must use $\Delta(t)$ instead. Clearly, if $\Delta(t)$ is nonzero but $z(t)$ is zero, the relative error $\delta(t)$ is infinite. That's an artifact of this example, and tells us something about the model.

Equation (4.14) is quite interesting. It gives the minimum possible backward error for the piece of the skeleton (t_n, t_y) to (t_{n+1}, y_{n+1}) . If we solve $\dot{z} = -\sqrt{z}(1 + \delta_n)$ on this interval, we have the “best possible” interpolant to this piece of the skeleton. The solution here is just as easy as the reference solution: $z(t) = (\sqrt{y_n} - (1 + \delta_n)(t - t_n)/2)^2$. Most importantly this analysis is independent of the method used to generate the skeleton. It applies to forward Euler, backward Euler, or indeed to any other method.

For both forward Euler and Euler backward one can use (4.14) to show that $\|\delta(t)\| = \mathcal{O}(\Delta t)$ as $\Delta t \rightarrow 0$. Together with the Gröbner–Aleksiev nonlinear variation of constants formula this establishes convergence of $z(t)$ to $y(t)$ with forward error (on $0 \leq t \leq 2\sqrt{y_0} - \varepsilon$ for every $\varepsilon > 0$, since we need G to be bounded) that goes to 0 as $\mathcal{O}(\Delta t)$, also. However, in practice one never takes $\Delta t \rightarrow 0$. One always uses a finite number of steps, and indeed one wants to take steps Δt as large as possible, for efficiency. This often means that we are outside the asymptotic ($\Delta t \rightarrow 0$) region.

A very important feature of (4.14) (and the similar equations that can be developed for other problems, even just using suboptimal interpolants) is that δ_n is computable. In this simple example, so is $z(t)$; in general that's harder, though. This gives a good retrospective diagnostic: after doing the computation, we can assess δ_n to see if our results are good enough. How does it work for these examples?

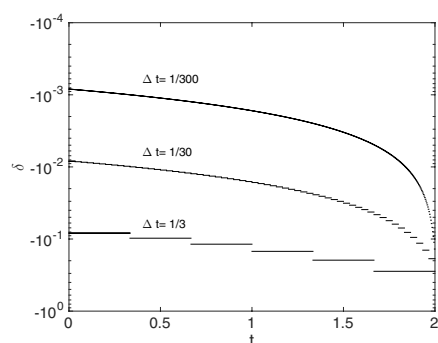


Fig. 4.1 Optimal backward error for the implicit Euler method for Torricelli's law, $\dot{y} = -\sqrt{y}$, $y(0) = 1$. We plot $\delta_n = 2 \cdot (\sqrt{y_n} - \sqrt{y_{n+1}})/\Delta t - 1$ for solutions using three different fixed time steps, $\Delta t = 1/3$, $\Delta t = 1/30$, $\Delta t = 1/300$. We see convergence away from the singularity when the bucket empties at $t = 2$. We also see the error become large near $t = 2$. Notice that $\delta_n < 0$ on this semilog scale.

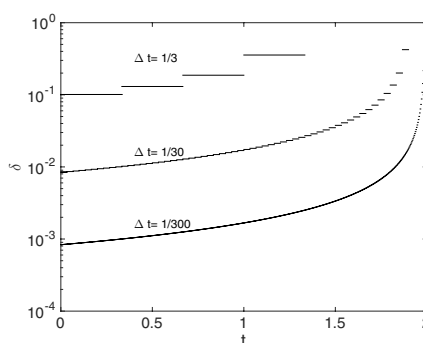


Fig. 4.2 Optimal backward error for the explicit Euler method for Torricelli's law, $\dot{y} = -\sqrt{y}$, $y(0) = 1$. We plot $\delta_n = 2 \cdot (\sqrt{y_n} - \sqrt{y_{n+1}})/\Delta t - 1$ for solutions using three different fixed time steps, $\Delta t = 1/3$, $\Delta t = 1/30$, $\Delta t = 1/300$. We see convergence away from the singularity when the bucket empties at $t = 2$. We also see the error become large near $t = 2$. Notice that for $\Delta t = 1/3$ the last two δ_n are blank because $y_5 < 0$.

Consider Figures 4.1 and 4.2. We see that for *both* methods the relative error rises as $y_n \rightarrow 0$. This is not really a surprise. What's nice is that for, say, $\Delta t = 1/30$, we see that the equation we've solved is, on $0 \leq t \leq 1.9$, no more than 5% different to the one we wanted to solve. This simple tool can be used for many equations.

Still, somehow, this “optimal backward error” might not satisfy someone who thinks of the differential equation as being “smooth”—here we have obvious pulses, changing with every time step. There's another, older, technique that can help, which we take up in section 6. But first, we look at optimal residuals for MATLAB's solution.

5. Optimal Backward Error for MATLAB's Solutions. The technique of optimal backward error can be used to analyze the skeleton produced by any numerical method. In Figures 5.1–5.6 we see the results from analyzing ode45, ode15s, and ode113, the three most popular MATLAB codes [15], each applied to the leaky bucket. We use the default tolerances for each method. In Figure 5.1 we plot not the

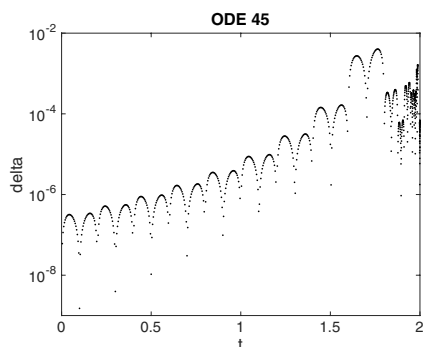


Fig. 5.1 We used `ode45` to solve $\dot{y} = -\sqrt{y}$, $y(0) = 1$, on $0 \leq t \leq 2$ using the default tolerances. Afterwards, we sampled the supplied interpolant and its derivative by using `deval` and plotted the relative residual $R(t) = (\dot{z} + \sqrt{z})/\sqrt{z}$ on $0 \leq t \leq 2$. The code to do this is available in the supplementary material. We see that the residual is initially small, but grows large near the singularity at the end, indicating either a problem with the skeleton (given by a 4-5 Runge-Kutta pair) or with the supplied interpolant, or both.

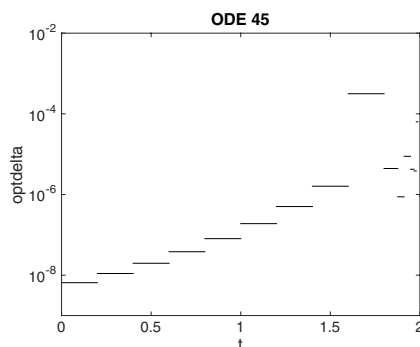


Fig. 5.2 We took the skeleton from the `ode45` solution as in Figure 5.1 and computed the optimal backward error using (4.14) on each step. Notice that here the optimum is about 10 times smaller than the results in Figure 5.1. This suggests that the supplied interpolant for `ode45` is not as accurate (indeed, it's only fourth order, but it suffices for many purposes and is quite convenient). The relative residual is never larger than 10^{-3} , which provides reassurance that the computed solution is valid.

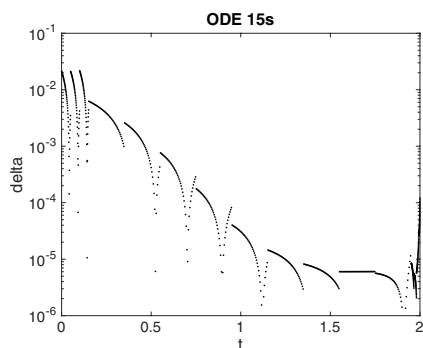


Fig. 5.3 We used `ode15s`, which is suited to stiff problems, to solve $\dot{y} = -\sqrt{y}$, $y(0) = 1$, using its default tolerances. As for Figure 5.1 and `ode45`, we sampled the supplied interpolant and its derivative and computed the relative residual $R(t) = (\dot{z} + \sqrt{z})/\sqrt{z}$. We see here that the residual is initially large-ish, as is suitable for the default tolerances which are quite loose, but improves as t gets larger. This is evidence that the problem gets stiff as z gets small, because `ode15s` performs better the stiffer the problem is.

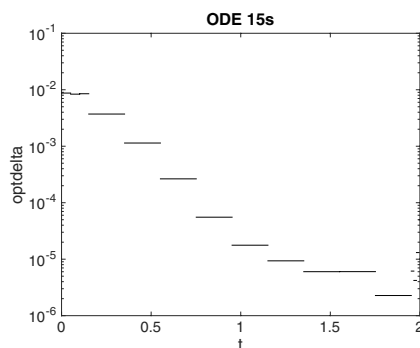


Fig. 5.4 We took the skeleton for the `ode15s` solution and computed the optimal backward error using (4.14) on each step. Notice here that the relative residual of the supplied interpolant plotted in Figure 5.3 is not much larger than the optimal, plotted here. This provides evidence that the supplied interpolant is high quality. That the backward error is small provides evidence of the accuracy of the solution.

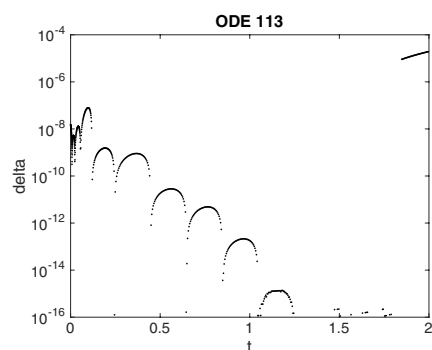


Fig. 5.5 We used `ode113` to solve $\dot{y} = -\sqrt{y}$, $y(0) = 1$, on $0 \leq t \leq 2$ using its default tolerances. Afterwards, we sampled the supplied interpolant and its derivative using `deval` and plotted the relative residual $R(t) = (\dot{z} + \sqrt{z})/\sqrt{z}$. We see that `ode113` is the most accurate method tried so far. We believe the extreme accuracy shown is an artifact of the simplicity of this problem: indeed a second order Taylor series method finds the reference solution exactly. The relative error on the last step is large, but this is only because the solution is so small.

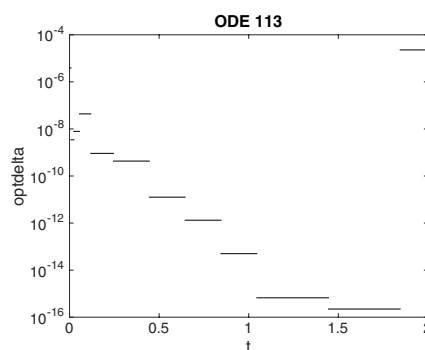


Fig. 5.6 We took the skeleton from the `ode113` solution and computed the optimal backward error using (4.14) on each step. The optimal backward error is hardly smaller than that from the supplied interpolants, showing that the supplied interpolants are of high quality. The optimal error on the final step is again large, relatively, but this is only because the solution is so small.

optimal backward error, but rather the backward error using the supplied piecewise polynomial interpolant. We see that it's fairly small, though not as small as the optimal, which is plotted in Figure 5.2. Similarly, Figure 5.3 plots the residual using the supplied interpolant for `ode15s`, while Figure 5.4 shows the optimal residual. Finally, Figures 5.5 and 5.6 show that for this problem `ode113` is very accurate indeed; even the supplied interpolant has only a tiny residual.

Remark 5.1. To do the computations we find it convenient to use the `RefineMesh` code from [4]. The code is included in the supplementary material.

In cases where the quality of the supplied interpolant is in doubt, computing the optimal backward error can supply useful reassurance.

6. The Method of Modified Equations. The method of modified equations has been studied at least since [10] and [18]. It is also discussed in [4], [9], and [11]. The idea is to start with the iteration defining the numerical method, say,

$$(6.1) \quad y_{n+1} = y_n - \Delta t \sqrt{y_n},$$

and conceptually interpolate with $z(t)$ so that $z(t_n) = y_n$ and $z(t_{n+1}) = y_{n+1}$. Now we try to find $z(t)$ so that the recurrence holds for all t with this fixed Δt . This converts the iteration to a functional equation:

$$(6.2) \quad z(t + \Delta t) = z(t) - \Delta t \sqrt{z(t)}.$$

We then use Taylor series, on the left-hand side, to convert this first to a *singularly* perturbed differential equation,

$$(6.3) \quad \begin{aligned} z(t) + \dot{z}(t)\Delta t + \frac{1}{2}\ddot{z}(t)\Delta t^2 + \cdots \\ = z(t) - \Delta t\sqrt{z(t)}, \end{aligned}$$

and then truncate and use repeated differentiation and substitution to convert this to a *regularly* perturbed differential equation. Just briefly, subtracting $z(t)$ from each side and dividing by Δt , we get

$$(6.4) \quad \begin{aligned} \dot{z}(t) + \frac{1}{2}\Delta t \ddot{z}(t) + \frac{1}{6}(\Delta t)^2 \dddot{z}(t) + \mathcal{O}(\Delta t^3) \\ = -\sqrt{z(t)}. \end{aligned}$$

Differentiating twice we get expressions for \ddot{z} and \dddot{z} that can be combined and manipulated (see supplementary material) to get

$$(6.5) \quad \left(-\frac{1}{\sqrt{z}} + \frac{\Delta t}{4z} + \mathcal{O}(\Delta t^3)\right) \frac{dz}{dt} = 1,$$

which has solution, expressed implicitly as an equation for the unknown $z(t)$, up to $\mathcal{O}(\Delta t^3)$

$$(6.6) \quad 2\sqrt{y_n} - 2\sqrt{z(t)} + \frac{\Delta t}{4} \ln\left(\frac{z(t)}{y_n}\right) = t - t_n.$$

Because $y_{n+1} = y_n - \Delta t\sqrt{y_{n+1}}$ can be written as $y_n = y_{n+1} + \Delta t\sqrt{y_{n+1}}$, we see that the modified functional equation for Euler backward is

$$(6.7) \quad y(t - \Delta t) = y(t) - (-\Delta t)\sqrt{y(t)},$$

which is the same equation as for forward Euler but with reversed sign of Δt . Therefore, one modified equation for Euler backward is

$$(6.8) \quad \left(-\frac{1}{\sqrt{z}} - \frac{\Delta t}{4z} + \mathcal{O}(\Delta t^3)\right) \frac{dz}{dt} = 1,$$

so

$$(6.9) \quad 2\sqrt{y_n} - 2\sqrt{z(t)} - \frac{\Delta t}{4} \ln\left(\frac{z(t)}{y_n}\right) = t - t_n.$$

We remark that modified equations are not unique. The difference in sign in the logarithm term has quite profound implications; see Figure 6.1. Notice that even for $\Delta t = 1/3$ these modified equations explain the numerical solution very well: it so happens that the $\mathcal{O}(\Delta t^3)$ term is $-\Delta t^3/192z^2$, leading to $(\Delta t^3)/192z$ on integration, which is very small as long as z is not near 0.

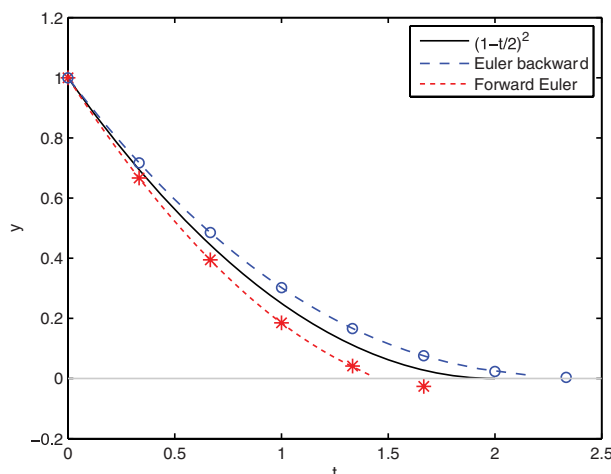


Fig. 6.1 Modified equations explaining forward and backward Euler for $\dot{y} = -\sqrt{y}$. The skeletons of the numerical solutions are shown with symbols \circ and $*$. They are nearly interpolated by the solutions of (6.11) and (6.10). $\Delta t = 1/3$.

In other words, forward Euler more nearly solves

$$(6.10) \quad \left(-\frac{1}{\sqrt{z}} + \frac{\Delta t}{4z} \right) \frac{dz}{dt} = 1,$$

while Euler backward more nearly solves

$$(6.11) \quad \left(-\frac{1}{\sqrt{z}} - \frac{\Delta t}{4z} \right) \frac{dz}{dt} = 1.$$

7. Returning to the Physics. Equation (1.2) sets the loss of potential energy equal to the gain in kinetic energy, with no dissipation. What if we added a crude model of dissipation, say, replacing $\frac{1}{2}(\rho A \Delta h)v^2$ by $\frac{1}{2}(\rho A \Delta h)v^2 - \varepsilon(\rho A \Delta h)v$, where ε is a small parameter, where the mass is present, times velocity⁸ as is common in friction terms? This model is hopelessly crude—friction and energy dissipation are not mere linear functions of momentum—but at least this model is simple and is somewhat like models that have been used before. Then setting $\rho A \Delta h \cdot gh$ equal to this gives

$$(7.1) \quad gh = \frac{1}{2}v^2 - \varepsilon v.$$

This has roots

$$(7.2) \quad v = +\varepsilon \pm \sqrt{\varepsilon^2 + 2gh}$$

and again we discard the positive root because h is decreasing; thus

$$(7.3) \quad \frac{A}{a} \frac{dh}{d\tau} = +\frac{\varepsilon}{H} - \sqrt{\frac{2g}{H}} \sqrt{\left(\frac{\varepsilon^2}{2gh} \right) + y}.$$

⁸Remember, $v < 0$. Thus the total change in potential energy $= 1/2mv^2 - \varepsilon mv$ means smaller magnitude velocity for a given potential difference, as is correct with dissipation present.

To clean things up put $\mu^2 = \frac{4\varepsilon^2}{2gH}$ or $\varepsilon = \sqrt{2gH}\mu$, so

$$(7.4) \quad \therefore \quad \frac{A}{a} \frac{dy}{d\tau} = +\sqrt{\frac{2g}{H}} - \sqrt{\frac{2g}{H}} \sqrt{\mu^2 + y}.$$

Putting $t = \frac{a}{A} \sqrt{\frac{2g}{H}} \tau$ as before makes

$$(7.5) \quad \frac{dy}{dt} = -(-\mu + \sqrt{\mu^2 + y}) = -\frac{y}{\mu + \sqrt{\mu^2 + y}},$$

or

$$(7.6) \quad -\left(\frac{\mu}{y} + \frac{\sqrt{\mu^2 + y}}{y}\right) \frac{dy}{dt} = 1.$$

Comparing this to the Euler backward modified equation (6.11) we see these will agree to $\mathcal{O}(\Delta t^2)$ if the dissipation $\mu \doteq \Delta t/4$. This is, Euler backward has introduced dissipation of about $\mu = \Delta t/4$ into the problem. This corresponds to $\varepsilon = \sqrt{2gH} \cdot \Delta t/4$ in the original equation.

What of forward Euler? Quite clearly its effect is qualitatively like negative dissipation—a kind of energy injection, if you will. The bucket might well empty faster if you shake it a bit!

This is the principal benefit of modified equations: they help to explain what numerical methods do, in terms that the modeler finds easier to understand.

8. Combining Optimal Backward Error and Modified Equations. One difficulty with modified equations is that they are truncated; another is that the series is not usually convergent, and is valid only in the asymptotic limit as $\Delta t \rightarrow 0$. Indeed, taking the next few terms here gives

$$\left(-\frac{1}{\sqrt{z}} - \frac{\Delta t}{4z} + \frac{(\Delta t^3)}{192z^2} - \frac{(\Delta t^4)}{384z^{5/2}} + \frac{(\Delta t^5)}{1920z^3} + \mathcal{O}(\Delta t^6)\right) \frac{dz}{dt} = 1.$$

(One might be tempted to conjecture that all the subsequent terms are simple reciprocal integers, but computing one more term shows that they're not.) We can see that if $z = \mathcal{O}(\Delta t^\beta)$ for $\beta < 1$, then these terms decrease, but at $z = \mathcal{O}(\Delta t^2)$ the terms actually increase; so we need to “close the loop” somehow. So we see that once z is as small as $\mathcal{O}(\Delta t^2)$, the modified equation is likely to be divergent if we take an infinite number of terms. We can argue that we should choose some fixed number of terms in a modified equation, e.g., (6.11),

$$(8.1) \quad \left(-\frac{1}{\sqrt{z}} - \frac{\Delta t}{4z}\right) \frac{dz}{dt} = 1,$$

and then ask for the optimal backward error $\hat{\delta}_n$ in *this* equation:

$$(8.2) \quad \left(-\frac{1}{\sqrt{z}} - \frac{\Delta t}{4z}\right) \frac{dz}{dt} = 1 + \hat{\delta}_n(t),$$

where $z(t_n) = y_n$ and $z(t_{n+1})$. Then this equation will completely explain the numerics. One can show that $\|\hat{\delta}_n\| = \mathcal{O}(\Delta t^3)$ as $\Delta t \rightarrow 0$. (By lucky chance the $\mathcal{O}(\Delta t^2)$

term is zero.) But even if Δt is “large,” we will have the best possible backward error. Integrating again,

$$(8.3) \quad \int_{\tau=t_n}^{t_{n+1}} \left(-\frac{1}{\sqrt{z(\tau)}} - \frac{\Delta t}{4z(\tau)} \right) \dot{z}(\tau) d\tau - \int_{\tau=t_n}^{t_{n+1}} d\tau = \int_{\tau=t_n}^{t_{n+1}} \hat{\delta}_n(\tau) d\tau,$$

which constrains $\delta(\tau)$ as before, because by the triangle inequality

$$(8.4) \quad \left| 2\sqrt{y_n} - 2\sqrt{y_{n+1}} - \frac{\Delta t}{4} \ln \left(\frac{y_n}{y_{n+1}} \right) - \Delta t \right| = \left| \int_{\tau=t_n}^{t_{n+1}} \delta(\tau) d\tau \right| \leq \int_{\tau=t_n}^{t_{n+1}} 1 d\tau \cdot \|\hat{\delta}_n\|_\infty$$

$$(8.5) \quad = \Delta t \|\hat{\delta}_n\|_\infty,$$

so

$$(8.6) \quad \|\hat{\delta}_n\|_\infty \geq \left| \frac{2}{\Delta t} (\sqrt{y_n} - \sqrt{y_{n+1}}) - \frac{1}{4} \ln \left(\frac{y_n}{y_{n+1}} \right) - 1 \right|$$

(we always take $\Delta t > 0$ here). This lower bound is attained if $\hat{\delta}_n(\tau)$ is constant, so

$$(8.7) \quad \begin{aligned} \hat{\delta}_n &= \left[\int_{t_n}^{t_{n+1}} \left(-\frac{1}{\sqrt{z}} - \frac{\Delta t}{4z} \right) \frac{dz}{dt} d\tau - \Delta t \right] \cdot \frac{1}{\Delta t} \\ &= \frac{2\sqrt{y_n} - 2\sqrt{y_{n+1}}}{\Delta t} - \frac{1}{4} \ln \left(\frac{y_{n+1}}{y_n} \right) - 1, \end{aligned}$$

which is very similar to (2.1). We see that $z(t)$ *exactly* solves

$$(8.8) \quad \left(-\frac{1}{\sqrt{z}} - \frac{\Delta t}{4z} \right) \frac{dz}{dt} = 1 + \hat{\delta}_n$$

on $t_n \leq t \leq t_{n+1}$. We also see that the optimal backward error here should be considerably smaller than if the term $\Delta t \ln(y_{n+1}/y_n)/4$ were not included. We also see that we need not worry about convergence: we have a complete model of the error that Euler backward introduces. We see finally in Figures 8.1 and 8.2 that it also gets (relatively) large as $z \rightarrow 0$, which is still not surprising.

We personally find this combined method of assessment of validity of numerical solutions very satisfactory. Modified equations by themselves make some people uneasy because of the asymptotics. Using optimal backward error as well removes that uneasiness.

9. Concluding Remarks. The main purpose of this paper was to use an understandable physical setting to explore the interpretation of numerical errors as modeling errors. The idea is surprisingly easy to get across to students, and it has happened more than once that a smiling student has asked during an exam, “Is it ok if I give the exact answer to a slightly different question?”

Of course, things become more complicated with larger models, and an especially serious issue is whether or not the $r(t)$ and $\delta(t)$ are physically reasonable perturba-

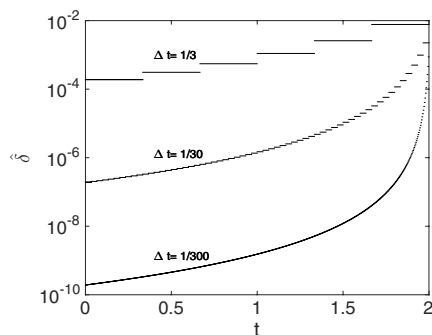


Fig. 8.1 Optimal backward error for the modified equation (6.10) that explains the bulk of the numerical effects of the implicit Euler method applied to Torricelli's law (3.4). Even for $\Delta t = 1/3$ we see that the residual numerical error is between 10^{-4} and 10^{-2} , showing that the bulk of the backward error shown in Figure 4.1 is explained by the term $-\frac{\Delta t}{4z} \cdot \frac{dz}{dt}$, incidentally illustrating that the implicit Euler method is first order.

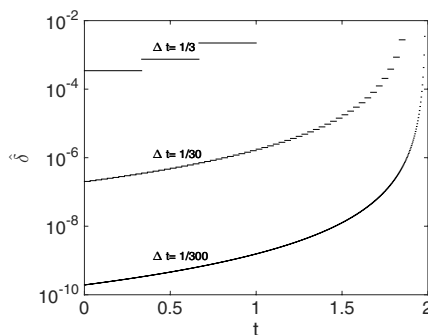


Fig. 8.2 Optimal backward error for the modified equation (6.7) that explains the bulk of the numerical effects of the forward Euler method applied to Torricelli's law (2.2). Even for $\Delta t = 1/3$ we see that the residual numerical error is much smaller than that in Figure 4.2, showing that the bulk of the backward error there is explained by the term $\frac{\Delta t}{4z} \cdot \frac{dz}{dt}$. Incidentally, this illustrates that Euler's method has first order.

tions; in other words, does the numerical method respect the structure of the model equations? This leads naturally to symplectic integration and other geometric integration methods [9], [11].

A relatively minor issue is the use of adaptive stepsize control. We originally had a section on this aspect for this paper but cut it for space reasons. Using the optimal backward error is straightforward in that context, and modified equations can be made to work, too (see [10]).

We have shown two ways to interpret what Euler's method does to this equation. Both ways work for other numerical methods, too: the optimal backward error is still given by (4.14), and modified equations can be constructed for any one-step method. We remark that some, but not all, second-order methods (second-order Taylor series, implicit trapezoidal rule, and implicit midpoint rule, but not RK2) actually give $r(t) \equiv 0$. That is, the numerical method gives the exact reference solution for Torricelli's law.

The Torricelli's law model, $\dot{y} = -\sqrt{y}$, is only one of many possible models of a leaky bucket. Another more complicated model is introduced in [2]. We critique that paper in [5].

In addition to the archive of supplementary material, our scripts and Maple worksheets can be found at www.apmaths.uwo.ca/~rcorless/Bucket.

Acknowledgments. J. M. Sanz-Serna pointed out that Euler knew that the explicit Euler method could fail and that Euler had a method for dealing with the kinds of singularity encountered here. His method amounts to using the reference solution for this example. Nic Fillion helped with the preparation of the figures and some of the text. Samantha Brennan's puppy, Cheddar, thought an early draft of the paper was chewy, palatable, and fun.

REFERENCES

- [1] G. BIRKHOFF AND G.-C. ROTA, *Ordinary Differential Equations*, Wiley, New York, 1989. (Cited on p. 777)
- [2] M. BLASONE, F. DELL'ANNO, R. D. LUCA, O. FAELLA, O. FIORE, AND A. SAGGESE, *Discharge time of a cylindrical leaking bucket*, Eur. J. Phys., 36 (2015), 035017, doi:10.1088/0143-0807/36/3/035017. (Cited on p. 791)
- [3] J. BUTCHER, *Numerical Methods for Ordinary Differential Equations*, Wiley, New York, 1987. (Cited on p. 777)
- [4] R. M. CORLESS AND N. FILLION, *A Graduate Introduction to Numerical Methods: From the Viewpoint of Backward Error Analysis*, Springer, New York, 2013. (Cited on pp. 780, 786)
- [5] R. M. CORLESS AND J. JANKOWSKI, *Revisiting the discharge time of a cylindrical leaking bucket*, in preparation (2016). (Cited on p. 791)
- [6] R. M. CORLESS AND C. Y. KAYA, *Minimizing residuals in ODE integration using optimal control theory*, in preparation (2016). (Cited on p. 782)
- [7] R. M. CORLESS, C. Y. KAYA, AND R. H. MOIR, *Optimal backward error and the Dahlquist test problem*, in preparation (2016). (Cited on p. 782)
- [8] W. H. ENRIGHT AND W. B. HAYES, *Robust and reliable defect control for Runge-Kutta methods*, ACM Trans. Math. Software, 33 (2007), 1206041. (Cited on p. 781)
- [9] D. F. GRIFFITHS AND D. J. HIGHAM, *Numerical Methods for Ordinary Differential Equations*, Springer, New York, 2010. (Cited on pp. 786, 791)
- [10] D. F. GRIFFITHS AND J. M. SANZ-SERNA, *On the scope of the method of modified equations*, SIAM J. Sci. Statist. Comput., 7 (1986), pp. 994–1008, doi:10.1137/0907067. (Cited on pp. 786, 791)
- [11] E. HAIRER, C. LUBICH, AND G. WANNER, *Geometric Numerical Integration: Structure-Preserving Algorithms for Ordinary Differential Equations*, Springer Series in Computational Mathematics, Springer, New York, 2006. (Cited on pp. 786, 791)
- [12] E. HAIRER, S. NØRSETT, AND G. WANNER, *Solving Ordinary Differential Equations I: Nonstiff Problems*, Springer, New York, 1993, <https://books.google.ca/books?id=F93u7VcSRyYC>. (Cited on pp. 777, 781)
- [13] J. H. HUBBARD AND B. H. WEST, *Differential Equations: A Dynamical Systems Approach. Part I. Ordinary Differential Equations*, Texts Appl. Math., Springer, New York, 1997. (Cited on pp. 775, 777, 782)
- [14] E. HUBERT, *Etude algébrique et algorithmique des singularités des équations différentielles implicites*, Ph.D. thesis, Institut National Polytechnique de Grenoble, 1997. (Cited on p. 779)
- [15] L. F. SHAMPINE AND M. W. REICHEL, *The MATLAB ODE suite*, SIAM J. Sci. Comput., 18 (1997), pp. 1–22, doi:10.1137/S1064827594276424. (Cited on p. 784)
- [16] G. SÖDERLIND, L. JAY, AND M. CALVO, *Stiffness 1952–2012: Sixty years in search of a definition*, BIT, 55 (2014), pp. 531–558. (Cited on p. 779)
- [17] L. N. TREFETHEN, *Numerical analysis*, in *The Princeton Companion to Mathematics*, T. Gowers, ed., 2010, pp. 604–615. (Cited on p. 780)
- [18] R. F. WARMING AND B. J. HYETT, *The modified equation approach to the stability and accuracy analysis of finite-difference methods*, J. Comput. Phys., 14 (1974), pp. 159–179. (Cited on p. 786)