Numerical Analysis – Lecture 14¹

Implementation 3.37 To implement the *Milne device*, we work with a *pair* of multistep methods of the same order, one explicit (*predictor*) and the other implicit (*corrector*), e.g.

Predictor:
$$\boldsymbol{y}_{n+2} = \boldsymbol{y}_{n+1} + h[\frac{5}{12}\boldsymbol{f}(t_{n-1},\boldsymbol{y}_{n-1}) - \frac{4}{3}\boldsymbol{f}(t_n,\boldsymbol{y}_n) + \frac{23}{12}\boldsymbol{f}(t_{n+1},\boldsymbol{y}_{n+1})],$$

Corrector: $\boldsymbol{y}_{n+2} = \boldsymbol{y}_{n+1} + h[-\frac{1}{12}\boldsymbol{f}(t_n,\boldsymbol{y}_n) + \frac{2}{3}\boldsymbol{f}(t_{n+1},\boldsymbol{y}_{n+1}) + \frac{5}{12}\boldsymbol{f}(t_{n+2},\boldsymbol{y}_{n+2})]$

the third-order Adams–Bashforth and Adams–Moulton methods respectively.

The predictor is employed not just to estimate the error of the corrector, but also to provide *an initial guess in the solution of the implicit corrector equations*. Typically, for nonstiff equations, we iterate correction equations at most twice, while stiff equations require *iteration to convergence*, otherwise the typically superior stability features of the corrector are lost.

Let TOL > 0 be a user-specified *tolerance*: the maximal error allowed in approximating the ODE. Having completed a single step and estimated the error, there are three possibilities:

- (a) $\frac{1}{10}$ TOL $\leq || \operatorname{error} || \leq \text{TOL}$, say: Accept the step, continue to t_{n+2} with the same step size.
- (b) $\| \operatorname{error} \| < \frac{1}{10}$ TOL, say: Accept the step and increase the step length;
- (c) $\| \operatorname{error} \| > \operatorname{TOL}$: Reject the step, recommence integration from t_n with smaller h.

Amending step size can be done easily with polynomial interpolation, although this means that we need to store past values well in excess of what is necessary for simple implementation of both multistep methods.

Improvement 3.38 Let e be our estimate of *local* error. Then e/h is our estimate for the global error in an interval of unit length. It is usual to require the latter quantity not to exceed TOL since good implementations of numerical ODEs should monitor the accumulation of *global* error. This is called *error estimation per unit step*.

Technique 3.39 (*Embedded Runge–Kutta methods*) The situation is more complicated with RK, since no single error constant determines local growth of the error. The approach of *embedded RK* requires, again, two (typically explicit) methods: an RK method of ν stages and order p, say, and another method, of $\nu + l$ stages, $l \ge 1$, and order p + 1, such that *the first* ν stages of both methods are identical. (This means that the cost of implementing the higher-order method is marginal, once we have computed the lower-order approximation.) For example, consider (and verify!)

$$\begin{aligned} & \mathbf{k}_{1} = \mathbf{f}(t_{n}, \mathbf{y}_{n}), \\ & \mathbf{k}_{2} = \mathbf{f}\left(t_{n} + \frac{1}{2}h, \mathbf{y}_{n} + \frac{1}{2}h\mathbf{k}_{1}\right), \\ & \mathbf{y}_{n+1}^{[1]} = \mathbf{y}_{n} + h\mathbf{k}_{2} & \implies \text{ order } 2, \\ & \mathbf{k}_{3} = \mathbf{f}(t_{n} + h, \mathbf{y}_{n} - h\mathbf{k}_{1} + 2h\mathbf{k}_{2}), \\ & \mathbf{y}_{n+1}^{[2]} = \mathbf{y}_{n} + \frac{1}{6}h(\mathbf{k}_{1} + 4\mathbf{k}_{2} + \mathbf{k}_{3}) & \implies \text{ order } 3. \end{aligned}$$

We thus estimate $\boldsymbol{y}_{n+1}^{[1]} - \boldsymbol{y}(t_{n+1}) \approx \boldsymbol{y}_{n+1}^{[1]} - \boldsymbol{y}_{n+1}^{[2]}$. [It might look paradoxical, at least at first glance, but the only purpose of the higher-order method is to provide error control for the lower-order one!]

Technique 3.39 (*The Zadunaisky device*) Suppose that the ODE y' = f(t, y), $y(0) = y_0$, is solved by an arbitrary numerical method of order p and that we have stored (not necessarily equidistant) past solution

¹Please email all corrections and suggestions to these notes to A.Iserles@damtp.cam.ac.uk. All handouts are available on the WWW at the URL http://www.damtp.cam.ac.uk/user/na/PartII/Handouts.html.

values $\boldsymbol{y}_n, \boldsymbol{y}_{n-1}, \dots, \boldsymbol{y}_{n-p}$. We form an interpolating *p*th degree polynomial (with vector coefficients) \boldsymbol{d} such that $\boldsymbol{d}(t_{n-i}) = \boldsymbol{y}_{n-i}, i = 0, 1, \dots, p$, and consider the differential equation

$$z' = f(t, z) + d'(t) - f(t, d), \qquad z(t_n) = y_n.$$
 (3.13)

There are two important observations with regard to (3.13)

(1) Since $d(t) - y(t) = O(h^{p+1})$, the term d'(t) - f(t, d) is usually small (because $y'(t) - f(t, y(t)) \equiv 0$). Therefore, (3.13) is a small perturbation of the original ODE.

(2) The exact solution of (3.13) is known: z(t) = d(t).

Now, having produced y_{n+1} with our numerical method, we proceed to evaluate z_{n+1} as well, using exactly the same method and implementation details. We then evaluate the error in z_{n+1} , namely $z_{n+1} - d(t_{n+1})$, and use it as an estimate of the error in y_{n+1} .

Problem 3.40: Solving nonlinear algebraic systems We have already observed that the implementation of an implicit ODE method, whether multistep or RK, requires the solution of (in general, nonlinear) algebraic equations in each step. For example, for an *s*-step method, we need to solve in each step the algebraic system

$$\boldsymbol{y}_{n+s} = \sigma_s h \boldsymbol{f}(t_{n+s}, \boldsymbol{y}_{n+s}) + \boldsymbol{v}, \qquad (3.14)$$

where the vector v can be formed from past (hence known) solution values and their derivatives. The easiest approach is *functional iteration*

$$\boldsymbol{y}_{n+s}^{[j+1]} = \sigma_s h \boldsymbol{f}(t_{n+s}, \boldsymbol{y}_{n+s}^{[j]}) + \boldsymbol{v}, \qquad j = 0, 1, \dots,$$

where $y_{n+s}^{[0]}$ is typically provided by the predictor scheme. It is very effective for *nonstiff* equations but fails for *stiff ODEs*, since the convergence of this iterative scheme requires similar restriction on h as that we strive to avoid by choosing an implicit method in the first place!

If the ODE is stiff, we might prefer a Newton-Raphson method, namely

$$\boldsymbol{y}_{n+s}^{[j+1]} = \boldsymbol{y}_{n+s}^{[j]} - \left[I - \sigma_s h \frac{\partial \boldsymbol{f}(t_{n+s}, \boldsymbol{y}_{n+s}^{[j]})}{\partial \boldsymbol{y}}\right]^{-1} [\boldsymbol{y}_{n+s}^{[j]} - \sigma_s h \boldsymbol{f}(t_{n+s}, \boldsymbol{y}_{n+s}^{[j]}) - \boldsymbol{v}].$$

The justification of the above is as follows: suppose that $y_{n+s}^{[j]}$ is an approximation to the solution. We linearise (3.14) locally about $(t_{n+s}, y_{n+s}^{[j]})$,

$$\boldsymbol{y}_{n+s} - \sigma_s h \boldsymbol{f}(t_{n+s}, \boldsymbol{y}_{n+s}) - \boldsymbol{v} \approx [\boldsymbol{y}_{n+s}^{[j]} - \sigma_s h \boldsymbol{f}(t_{n+s}, \boldsymbol{y}_{n+s}^{[j]}) - \boldsymbol{v}] + \left[I - \sigma_s h \frac{\partial \boldsymbol{f}(t_{n+s}, \boldsymbol{y}_{n+s}^{[j]})}{\partial \boldsymbol{y}}\right] (\boldsymbol{y}_{n+s} - \boldsymbol{y}_{n+s}^{[j]})$$

and choose $oldsymbol{y}_{n+s}^{[j+1]}$ by equating the right-hand side to $oldsymbol{0}.$

The snag is that repeatedly evaluating and inverting (i.e. LU-factorizing) the Jacobian matrix in every iteration is *very* expensive. The remedy is to implement the *modified Newton–Raphson method*, namely

$$\boldsymbol{y}_{n+s}^{[j+1]} = \boldsymbol{y}_{n+s}^{[j]} - \left[I - \sigma_s h \frac{\partial \boldsymbol{f}(t_{n+s}, \boldsymbol{y}_{n+s}^{[0]})}{\partial \boldsymbol{y}}\right]^{-1} [\boldsymbol{y}_{n+s}^{[j]} - \sigma_s h \boldsymbol{f}(t_{n+s}, \boldsymbol{y}_{n+s}^{[j]}) - \boldsymbol{v}].$$
(3.15)

Thus, the Jacobian need be evaluated only *once* a step.

The only role the Jacobian matrix plays in (3.15) is to ensure convergence: its precise value makes no difference to the ultimate value of $\lim_{j\to\infty} y_{n+s}^{[j]}$. Therefore we might replace it with a finite-difference approximation, evaluate it once every several steps etc. Further note that the implementation of the method requires the LU factorization of $I - \sigma_s h J$, where J is our approximation of the Jacobian matrix, only once we change either h or J. Recalling that, once the LU factorization is known, the solution of linear systems is cheaper by an order of magnitude than the LU factorization itself, we deduce that the overall cost of this procedure is not excessive. For stiff equations it is definitely cheaper than using a minute step size with a 'bad' (e.g., explicit multistep or explicit RK) method.