## Numerical Instability

Consider the initial-value problem

(1)
$$\frac{dY}{dt} = -KY, \quad Y(0) = Y_0, \quad K > 0$$

which of course has the solution

(2)
$$Y(t) = Y_0 e^{-Kt}$$

Euler's method for this problem is
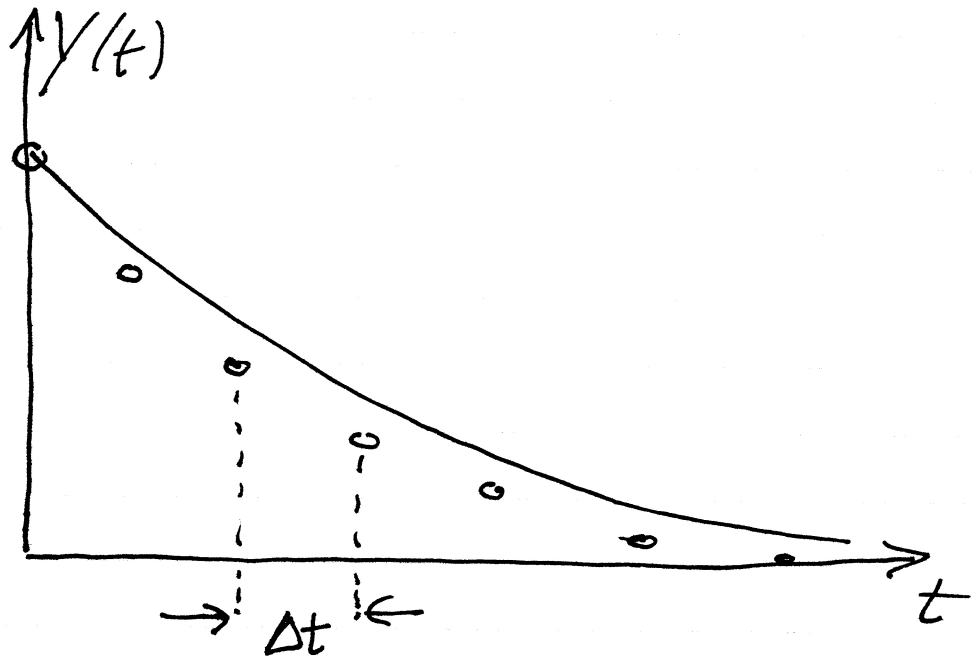
(3)
$$\frac{Y(t + \Delta t) - Y(t)}{\Delta t} = -KY(t)$$

and this implies

(4)
$$Y(t + \Delta t) = \left(1 - (\Delta t)K\right) Y(t)$$

So

(5)
$$Y(n\Delta t) = \left(1 - (\Delta t)K\right)^n Y(0)$$

If $(\Delta t)K < 1$, the solutions (2) and (5) are qualitatively the same in that both decay:



but if $(\Delta t)K > 1$ the solutions of the Euler scheme changes sign at every step, and if $(\Delta t)K > 2$ the solution actually grows in magnitude while still alternating sign.

This numerical <u>instability</u> when $\Delta t$ is too large has nothing to do with the original continuous problem; it is an artifact of the numerical method.

In the above simple problem, we would never want to choose $(\Delta t) K > 2$, since we need $(\Delta t) K \ll 1$ for accuracy of the computed solution, but what if $Y$ is a vector and $K$ is a matrix. (For simplicity, think of the case in which $K$ is a symmetric positive-definite matrix, so we

can do an eigenvector expansion of $Y(t)$ and each component behaves as above, but with $K$ replaced by an eigenvalue of $K$.)

Now there are multiple time scales, and we may not want to resolve them all, but if we use Euler's method, the fastest time scale will determine the time step.

The <u>backward-Euler</u> method for equation (1) is

(6)
$$\frac{Y(t) - Y(t - \Delta t)}{\Delta t} = -K Y(t)$$

Solving for $Y(t)$, we get

(7)
$$Y(t) = \frac{Y(t - \Delta t)}{1 + (\Delta t) K}$$

when $(\Delta t) K << 1$, this is essentially the same as (4), but when $(\Delta t) K$ is large, the two are very different! Here $\Delta t$ can be arbitrarily large, and the qualitative behavior is still correct.

Note that when $K$ is a matrix and $Y$ is a vector the implementation of the backward-Euler method requires the solution of a linear system at each step, since (7) becomes

(8) $$Y(t) = \left(I + (\Delta t) K\right)^{-1} Y(t - \Delta t)$$

Also, if the original problem was nonlinear

(9) $$\frac{dY}{dt} = f(Y)$$

Then we have a non-linear system to solve :

(10)
$$\frac{Y(t) - Y(t - \Delta t)}{\Delta t} = f(Y(t))$$

or

(11)
$$Y(t) - \Delta t \, f(Y(t)) = Y(t - \Delta t)$$

If $Y(t)$ is a vector, this is a nonlinear <u>system</u> of equations. One method of solution is Newton's method, in which we make a guess for $Y(t)$, linearize the problem around that guess by Taylor series, and then solve the linear problem to get the next guess.