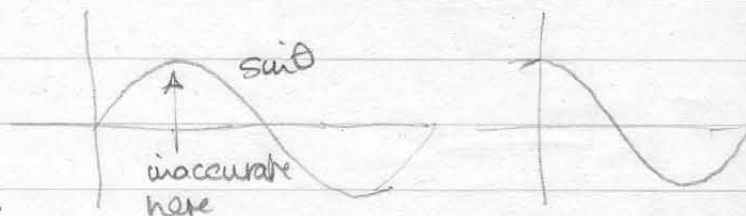


1. Plots
2. Access to Mathematica / computers.
3. Will use algebraic manip program, Will plot on cptr
3. How to solve quadratic?
4. Moral: avoid getting solution from difference of 2 large numbers.
5. Given $\sin \theta$ and $\cos \theta$, what is numerically accurate way to get θ ? \sin^{-1} ? \cos^{-1} ? \tan^{-1}

6. Why?



7. $\sinh \theta$ & $\cosh \theta$?

8. Given cubic eq $ax^3 + bx^2 + cx + d = 0$, how would you solve it?

not so useful

(a) Analytic solution?

(b) Newton-Raphson?

(c) Root finder?

← involves soln of quadratic
 fastest, esp if you need all roots. Need \tan^{-1} or \sinh^{-1} .
 ← next fastest. Need init solution. Does NR work for \mathbb{C} ?
 ← most generic.

9. Quartic?

(a) analytic - involves soln of cubic, then quadratic must polish real cubic root before completing.

10. General polynomial? All roots, including complex?
 NR recommend "Laguerre" iteration, See NR followed by forward deflation. See NR.
 Convergence is slowed if roots are multiple.

Solution of quadratic

$$ax^2 + 2bx + c = 0$$

$$x = \frac{-b \pm \sqrt{b^2 - ac}}{a}$$

Assume $ab \gg 0$ wlog
can get large cancellation.

$\frac{-b - \sqrt{b^2 - ac}}{a}$ is fine ($b > 0$)

but $\frac{-b + \sqrt{b^2 - ac}}{a}$ gives large cancellation
if $b^2 \gg ac$.

Alt:

$$a + \frac{2b}{x} + \frac{c}{x^2} = 0$$

$$\frac{1}{x} = \frac{-b \pm \sqrt{b^2 - ac}}{c}$$

$x = \frac{c}{-b \pm \sqrt{b^2 - ac}}$ is well-behaved
alt to

Or,

$$x = \frac{-b + \sqrt{b^2 - ac}}{a}$$

$$= \frac{(-b + \sqrt{b^2 - ac})(-b - \sqrt{b^2 - ac})}{a(-b - \sqrt{b^2 - ac})}$$

$$= \frac{b^2 - (b^2 - ac)}{a(-b - \sqrt{b^2 - ac})} = \frac{ac}{a(-b - \sqrt{b^2 - ac})}$$

$$= \frac{c}{-b - \sqrt{b^2 - ac}} \quad \checkmark$$

Extra]

Interpolation of functions

various parts of Chs 3, 4, 5 of Press.

Given a set of points, fit a ^{(reasonably) smooth} line through them.

Common types of fit:

- Polynomial
- Rational function
- Chebyshev polynomial
- Spline

never use it.

Why do you need fit?

- Interpolation
- Integration (quadrature)
- ODEs
- PDEs

v. important!
FoundationNeed to solve continuum
result in finite
number of chunksPolynomial fit

differentiable

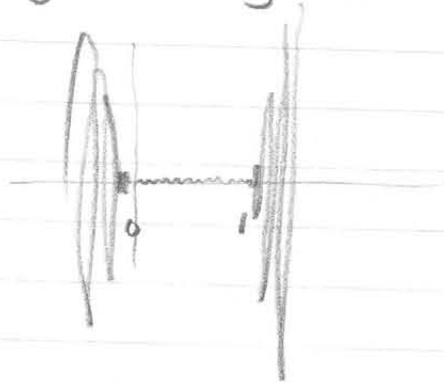
Justification: most functions 'of interest' are analytic:

$$f(z) = f(z_0) + \underset{\substack{\uparrow \\ z-z_0}}{\Delta z} f'(z_0) + \frac{\Delta z^2}{z!} f''(z_0) + \dots$$

As $|\Delta z| \rightarrow 0$, approxn by first few terms becomes increasing good.

BUT (i) Your function may not be well-behaved
eg. it may have singularity,
(ii) or maybe it's not smooth

(ii) \exists polynomials which are arbitrarily good approximations to $f(x) = 0$ in $x \in [0, 1]$, but go crazy outside $[0, 1]$.



Beware of extrapolation!

Lagrange formula (which you all know)

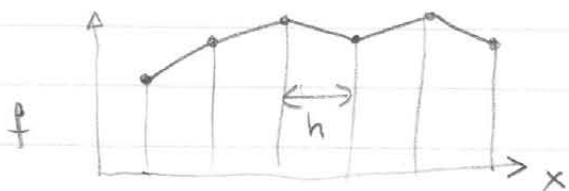
$$f_0 \frac{(x-x_1) \dots (x-x_N) f_0 + (x-x_0)(x-x_2) \dots (x-x_N) f_1 + \dots}{(x_0-x_1) \dots (x_0-x_N)} + \dots + \frac{(x-x_0) \dots (x-x_{N-1}) f_N}{(x_N-x_0) \dots (x_N-x_{N-1})}$$

is polynomial of degree N through $N+1$ points
 $f(x_n) = f_n$

Press gives Neville's algorithm = neat way of building up Lagrange formula recursively.

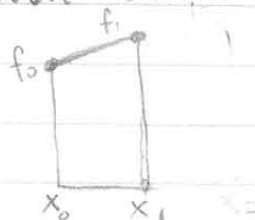
Application of polynomial fitting to integration
 Take equally spaced intervals h

1. Trapezium rule



Interpolate each neighboring pair with linear polynomial

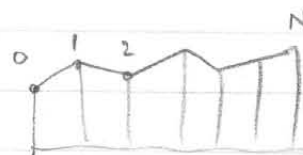
Want $\int_{x_0}^{x_1} f(x) dx$



$f = f_0(1-t) + f_1 t$ $t \equiv \frac{x-x_0}{x_1-x_0} \leftarrow h$

$$\begin{aligned} \int_{x_0}^{x_1} f(x) dx &= h \int_0^1 f(t) dt \\ &= h \int_0^1 (1-t)f_0 + t f_1 dt \\ &= h \left[\left(t - \frac{t^2}{2}\right) f_0 + \frac{t^2}{2} f_1 \right]_0^1 \\ &= h \left(\frac{f_0}{2} + \frac{f_1}{2} \right) \end{aligned}$$

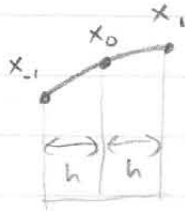
f is accurate to $O(h^2)$
 so $\int f dx$ is accurate to $O(h^3)$
 Sum N intervals: $h \sim 1/N$



weight on $f_n \rightarrow \frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \dots, \frac{1}{2}$ $Nh^3 \downarrow O(N^{-2})$

$$\int_{x_0}^{x_N} f dx = h \left(\frac{1}{2} f_0 + f_1 + \dots + f_{N-1} + \frac{1}{2} f_N \right) \sim O(N^{-2})$$

2. Simpson rule - some surprises

Fit quadratic to 3, ^{equally spaced} points

$$\int_{x_{-1}}^{x_1} f dx = h \int_{-1}^1 f dt \quad t = \frac{x - x_0}{x_1 - x_0} \leftarrow h$$

$$f = \frac{(t^2 - t)}{2} f_{-1} + (1 - t^2) f_0 + \frac{(t^2 + t)}{2} f_1 + O(t^3)$$

$$\int_{x_{-1}}^{x_1} f dx = h \left[\left(\frac{t^3}{6} - \frac{t^2}{4} \right) f_{-1} + \left(t - \frac{t^3}{3} \right) f_0 + \left(\frac{t^3}{6} + \frac{t^2}{4} \right) f_1 \right]_{-1}^1$$

$$= h \left(\frac{f_{-1}}{3} + \frac{4f_0}{3} + \frac{f_1}{3} \right)$$

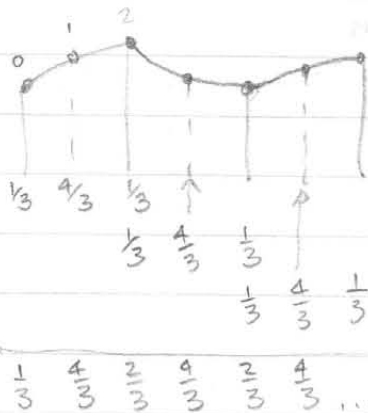
f is accurate to $O(h^3)$ so you'd think $\int_{x_{-1}}^{x_1} f dx$ accurate to $O(h^4)$

but actually if you add cubic term:

$$\int_{-1}^1 t^3 dt = \left[\frac{t^4}{4} \right]_{-1}^1 = 0 \quad \text{disappears}$$

 $\Rightarrow \int_{x_{-1}}^{x_1} f dx$ accurate to $O(h^5)$ surprise

Sum N intervals



$$\int_{x_0}^{x_N} f dx$$

$$= h \left(\frac{f_0}{3} + \frac{4f_1}{3} + \frac{2f_2}{3} + \frac{4f_3}{3} + \dots + \frac{4f_{N-1}}{3} + \frac{f_N}{3} \right)$$

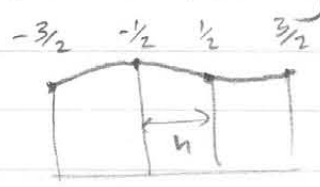
$$+ O(N^{-4})$$

$$\uparrow$$

$$Nh^5$$

surprise!
 Notice way weights alternate. Mysterious! Are alternating weights somehow better than equal weights? Alternation comes from asymmetry in way of treating intervals - some are left, some are right. Better idea to treat intervals on equal footing.

3. Press talks about Simpson's 3/8 rule, which comes from fitting cubic to 4 points:



$$\int_{x=-3/2}^{x=3/2} f dx = h \left(\frac{3}{8} f_{-3/2} + \frac{9}{8} f_{-1/2} + \frac{9}{8} f_{1/2} + \frac{3}{8} f_{3/2} \right)$$

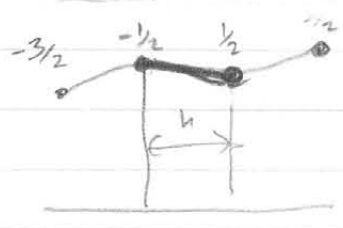
f is accurate to $O(h^4)$

So $\int_{x=-3/2}^{x=3/2} f dx$ is accurate to $O(h^5)$, like Simpson.

Press then meshes Simpson & Simpson 3/8 to get weights which are = 1 inside.

I prefer slightly different approach:

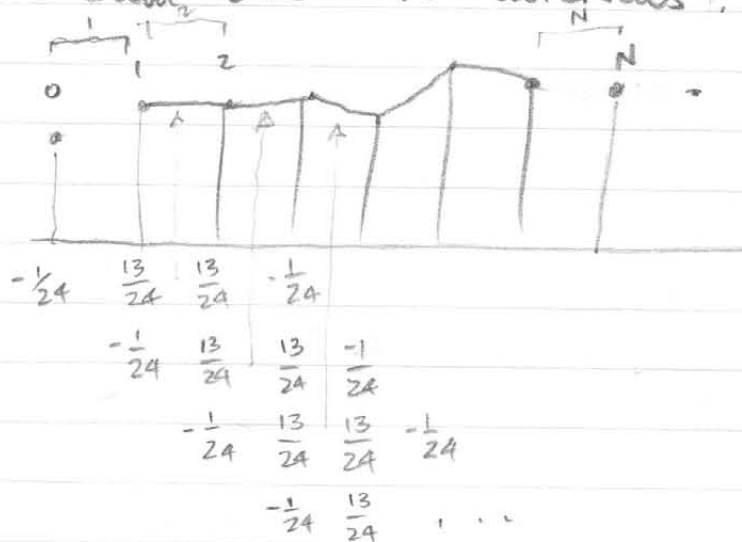
4. Like Simpson 3/8, fit cubic to 4 points:



but now look at integral only over central interval.

Result: $\int_{x=-1/2}^{x=1/2} f dx = h \left(-\frac{1}{24} f_{-3/2} + \frac{13}{24} f_{-1/2} + \frac{13}{24} f_{1/2} - \frac{1}{24} f_{3/2} \right) + O(h^5)$

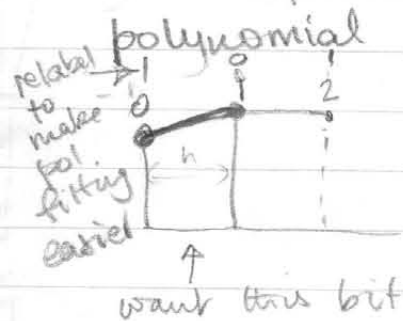
Now sum over N intervals:



Each interval is accurate to $O(N^{-5})$
 so sum is accurate to $O(N^{-4})$

$$-\frac{1}{24} \quad \frac{12}{24} \quad \frac{25}{24} \quad | \quad | \quad | \dots$$

But what to do about end intervals? Because only 2 end intervals, can take one order lower in polynomial fit and still get $O(N^{-4})$. (Press says you can't do this, but in fact you can).



So fit quadratic to first 3 points (same as Simpson), but take only first interval

$$\int_{x_{-1}}^{x_0} f dx = h \left(\frac{5}{12} f_{-1} + \frac{2}{3} f_0 - \frac{1}{12} f_1 \right)$$

$$-\frac{1}{24} \quad \frac{12}{24} \quad \frac{25}{24} \quad | \quad | \quad | \dots$$

+ endpoint contribution

$$\frac{10}{24} \quad \frac{16}{24} \quad \frac{-2}{24}$$

$$\frac{9}{24} \quad \frac{28}{24} \quad \frac{23}{24}$$

$$\Rightarrow \int_{x_0}^{x_N} f dx = h \left(\frac{3}{8} f_0 + \frac{7}{6} f_1 + \frac{23}{24} f_2 + f_3 + \dots + \frac{3}{8} f_N \right) + O(N^{-4})$$

Question: is it better to go to higher order or to reduce stepsize h ?

Answer: it depends. Going to higher order helps only if you know your function is sufficiently smooth. Δkh must certainly be smaller than radius of convergence.

As regards solution of Diff Eqs, must also worry about stability: higher order may mean more spurious solutions to Eqs, which may be disastrous.

Recommendation: a cubic is good compromise in most everyday work. Linear is best when function is not smooth. For a few problems - eg. integration of solar system for long time - higher order scheme may be best.

Rational function interpolation

$R(x) = \frac{P(x)}{Q(x)}$ \hookrightarrow polynomials is rational function.

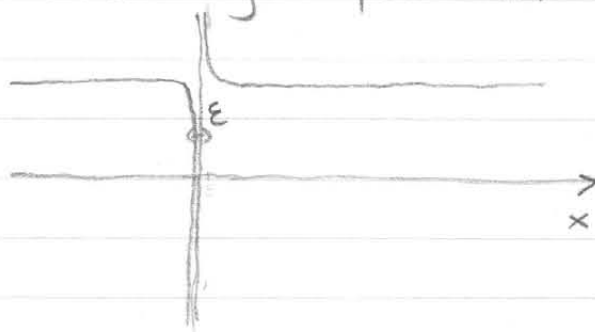
1. Rat. fn. is superior to polynomial in ability to model nearby poles. May work better for your function; may permit longer leaps.

Press clearly likes 'em... but as always "look before you leap"! Advantage

extrapolation to $x = \frac{1}{10}$ better behaved
fit over here

2. Beware of near cancelling factors:

eg. $\frac{x+\epsilon}{x}$
Small number



≈ 1 except within $\sim \epsilon$ of 0.



looks like great fit, but it ain't.

My Experience says this happens.

Disadvantage of rational functions as opposed to polynomials.

3. Press gives "Neville" type ^{recursive} algorithm for generating rat fn fit to set of points.

4. How to solve fit to rational fn?

Ans: rewrite as $RQ = P$, solve resulting linear eqn.

In case $n=1$

$$q_0 = 1 \quad (\text{of course})$$

$$r_1 q_1 = -r_2 \quad \Rightarrow \quad q_1 = -r_2 / r_1$$

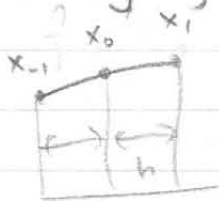
$$p_0 = r_0 q_0 = r_0$$

$$p_1 = r_1 q_0 + r_0 q_1 = r_1 + \frac{r_0 r_2}{r_1}$$

$$\Rightarrow \frac{P}{Q} = \frac{r_0 + (r_1 - \frac{r_0 r_2}{r_1})x}{1 - \frac{r_2}{r_1}x}$$

$$= \frac{r_0 r_1 + (r_1^2 - r_0 r_2)x}{r_1 - r_2 x} = r_0 + \frac{r_1^2 x}{r_1 - r_2 x}$$

Interpolation through 3 equally spaced points:
Previously got quadratic pol fit



$$t = \frac{x - x_0}{x_1 - x_0} = h$$

$$f = f_0 + \left(\frac{f_1}{2} - \frac{f_{-1}}{2} \right) t + \left(\frac{f_{-1}}{2} - f_0 + \frac{f_1}{2} \right) t^2$$

$$= r_0 + r_1 t + r_2 t^2$$

$n=1$ rat fit

$$= f_0 \left(\frac{f_1 - f_{-1}}{2} \right) + \left[\left(\frac{f_1 - f_{-1}}{2} \right)^2 - f_0 \left(\frac{f_{-1}}{2} - f_0 + \frac{f_1}{2} \right) \right] t$$

$$\frac{f_1 - f_{-1}}{2} - \left(\frac{f_{-1}}{2} - f_0 + \frac{f_1}{2} \right) t$$

Sort of yuk.

Spline

Disadvantage of polynomial or rational interpolation is that derivatives are discontinuous at interpolation points.



Aim of spline is to make derivs continuous.

Most popular is cubic spline.

Idea: fit cubic polynomial between ^{each} pair of points



$$f \approx p_0 + p_1 x + p_2 x^2 + p_3 x^3 = P(x)$$

Adjust 4 coefficients so

(a) $P(x)$ goes through $f(x)$ at $x = x_0$ & x_1 ,

(b) $\frac{dP}{dx}$ is continuous at x_0 & x_1 , with P of adjacent interval.

= 4 conditions, except at endpoints.

At ^{each of 2} endpoints, need to impose 1 more condition

eg. (c) $\frac{dP}{dx^2} = 0$ at endpoints

Requirement that $\frac{dP}{dx}$ is continuous across interp points

is non-local condition. End up needing to solve N linear equations, which look like $A_{ij} f_j'' = B_i$

- see Press for details.

A_{ij} is tridiagonal \rightarrow can be

solved in $O(N)$ steps = quick.

↑
functions of x_k, f_k

Press gives excellent discussion - if you need spline, read him.

Evaluation of functions

Summary of methods:

(1) Power series

- OK where series converges sufficiently strongly.

(2) Transformations:

eg. $\sin(x + 2\pi) = \sin x$

$$\Gamma(1+z) = z\Gamma(z)$$

$$G(\frac{1}{z}) = -G(z) + \frac{1}{2} \ln^2 z + \frac{\pi^2}{6}$$

(3) Recurrence relations:

transform a function into others which you have a better hope of evaluating

eg. Legendre polynomial

$$(l+1)P_{l+1}(x) = (2l+1)xP_l(x) - lP_{l-1}(x)$$

starting from $P_{-1} = 0, P_0 = 1$

Bessel function

$$J_{n+1}(x) = \frac{2n}{x} J_n(x) - J_{n-1}(x)$$

(4) Continued fraction

= infinite rational approximation. eg. incomplete $\times \frac{\Gamma n}{\beta}$

(5) Asymptotic series.

= series diverges, but first few terms give accurate result for sufficiently extreme values of argument.

We'll be finding each of these works best, depending on function, and value of argument of function.

$$(2') \text{ Incomplete } \beta\text{-fn } B_x(a,b) = \frac{x^a(1-x)^b}{a} + \frac{a+b}{a} B_x(a+1,b) = \int_0^x t^{a-1}(1-t)^{b-1} dt$$

$$= -\frac{x^a(1-x)^b}{b} + \frac{a+b}{a} B_x(a,b+1)$$

Note $B(a,b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$ for $a, b > 0$; extend by analytic cont.

Nth order Runge-Kutta method

Integrate $\frac{dy}{dx} = f(x, y)$

How does this differ from quadrature?



numerically from x to x + Δx.

1. Basic assumption:

Each y approximates Nth order polynomial in x.

2. Advantage: ^{simplicity} Depends only on evaluations within interval [x, x + Δx], not on history.

3. Disadvantage: Not fastest, as does not use history of previous evaluations.

Not most robust, nor most efficient.

Example: 2nd order RK

For simplicity take single y, but carries through to vector of y's.

Taylor series solution:

$$y(x + \Delta x) = y + \frac{dy}{dx} \Delta x + \frac{d^2y}{dx^2} \frac{\Delta x^2}{2!} + O(\Delta x^3)$$

$$\frac{df}{dx} = \frac{\partial f}{\partial x} + \frac{dy}{dx} \frac{\partial f}{\partial y}$$

so

$$y(x + \Delta x) = y + f \Delta x + \left(\frac{\partial f}{\partial x} + f \frac{\partial f}{\partial y} \right) \frac{\Delta x^2}{2!} + O(\Delta x^3)$$

Know all except this. How to estimate it?

1st order guess is

$$\Delta y = f \Delta x = f_0 \Delta x$$

$$\begin{aligned} \text{Try } f_1 &= f(x + \Delta x, y + \Delta y) \\ &= f + \frac{\partial f}{\partial x} \Delta x + \frac{\partial f}{\partial y} \Delta y \\ &= f + \left(\frac{\partial f}{\partial x} + f \frac{\partial f}{\partial y} \right) \Delta x \end{aligned}$$

Why OK to linear order?

Hey, this looks like desired. Yes!

Organise this as

$$\begin{cases} f_0 = f(x, y) \\ f_1 = f(x + \Delta x, y + f_0 \Delta x) \end{cases}$$

Then

$$y(x + \Delta x) = y + \frac{1}{2} (f_0 + f_1) \Delta x$$

General RK

$$y(x + \Delta x) = y + \Delta x \sum_{i=0}^s a_i f_i$$

$$f_0 = f(x, y)$$

$$f_1 = f(x + c_1 \Delta x, y + b_{10} f_0 \Delta x)$$

⋮

$$f_i = f(x + c_i \Delta x, y + \Delta x \sum_{j < i} b_{ij} f_j)$$

c_1	b_{10}		
c_2	b_{20}	b_{21}	
⋮			
c_s	b_{s0}	⋯	$b_{s,s-1}$
	a_0	⋯	a_s

4th order RK

The most common (by far?) scheme is following 4th order RK.

$f_0 = f(x, y)$	
$f_1 = f(x + \frac{1}{2} \Delta x, y + \frac{1}{2} f_0 \Delta x)$	
$f_2 = f(x + \frac{1}{2} \Delta x, y + \frac{1}{2} f_1 \Delta x)$	
$f_3 = f(x + \Delta x, y + f_2 \Delta x)$	
$y = y + (\frac{1}{6} f_0 + \frac{1}{3} f_1 + \frac{1}{3} f_2 + \frac{1}{6} f_3) \Delta x$	

Run rk4.nb

-Moral: mathematica is pretty dumb, but good at checking whether something = 0

Spline

gimp iras-100um-2k.rgb

carina_acs-1920x1200.pnm

← Nathan Smith

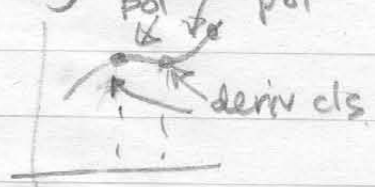
tools → color tools → curves

Fit to ps'

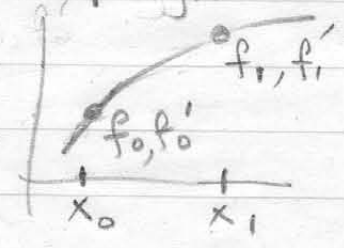
Spline fitting = closest approxn to "fitting w French curve"

A N^{th} spline is a fit through a set of points such that

1. The fit on each segment between each pair of adjacent points is a N^{th} order polynomial.
2. The derivative is continuous across segments.



Q. What is minimum order polyn that fits through $f(x_i)$ with specified values
 $f(x_0) = f_0, f'(x_0) = f'_0$
 $f(x_1) = f_1, f'(x_1) = f'_1$?



A. 3. ie, cubic. $ax^3 + bx^2 + cx + d$, 4 coeffs fits 4 things.

3. Can also require higher-order derivs to be order deriv at each, eg $f''(x_0) = f''_0, f''(x_1) = f''_1$ etc.

Q. Order polyn required to fit up to M^{th} deriv?

A. $N = 2M + 1$. \Rightarrow splines almost always odd order.

4. Most common (by far) spline is cubic, which fits values f_i and derivs f'_i .

5. What to do about derivs at endpoints?

For cubic order spline, choose 2 boundary conds. $(N-1)$

For cubic Most common choice: set $f''_0 = 0, f''_n = 0$, corresponding to extrapolation with straight line. But can also choose f'_0 & f'_n to be whatever.

Q: What do you need to know mathematically about spline fitting?

A: It's more complicated than polyn (or rat) fitting, but not ridiculously so.

For cubic spline through n points need to solve system of n linear equations. Eqs link i th pt to $i-1$ th and $i+1$ th pts, so $n \times n$ is tridiagonal.

Q: When might you use spline-fitting?

A: When you want fit that:

- (a) is robust (never fails)
- (b) has continuous derivs
- (c) gives decent (French curve) fit with few pts
- (d) extrapolates nicely outside prescribed interval.

Ex: fit atomic/molecular cross-sections from a modest number of tabulated pts, with prescribed asymptotic behavior.

Q: When not use spline-fitting?

A: In general numerical integrators,

- overhead of spline-solving too much
- more efficient to use polyn, more points.

Rational function fitting

A rational function $R(x)$ is a ratio of polynomials:

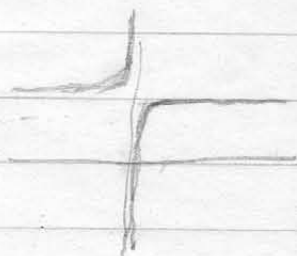
$$R(x) = \frac{P(x)}{Q(x)} \quad \left\{ \begin{array}{l} \text{order } m \text{ polyn} \\ \text{order } n \text{ polyn} \end{array} \right.$$

Q: When would you NOT use rat fn fitting? *

A: In any canned numerical routine.

Why?

$$\frac{x - \epsilon}{x} \rightarrow \begin{cases} 0 & \text{at } x = \epsilon \\ \pm \infty & \text{at } x = 0 \\ 1 & \text{at } x \rightarrow \pm \infty \end{cases}$$



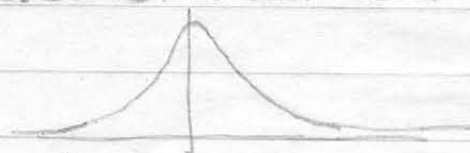
Q: When might you use rat fn fitting?

A: When you want to approximate some specific formula(e)

Q: What is advantage of rat fn fitting over polynomial fitting?

A: Can remain well-behaved over all x .

Eg. $y = \frac{1}{1+x^2}$



B: If there are poles, rats can do 'em.

Q: What do you need to know mathematically about rat fn fitting?

A: It's not that hard to do.

* Mathematica uses polyn interpolation to produce InterpolatingFunction objects.

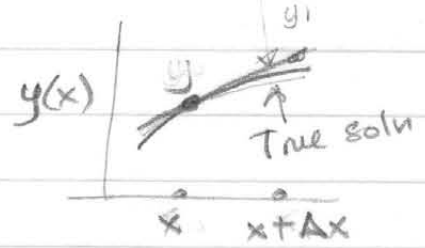
Diff Eq - numerical solution
 simplest possible scheme
 Take 1st order RK ("Euler")

- How does num solns compare to true?
- Is num soln stable?

Solve $\frac{dy}{dx} = f(x, y)$

by $f_0 = f(x_0, y_0)$

$y_1 = y(x + \Delta x) \approx y_0 + f_0 \Delta x$



Q: Does this come from polyn fitting? A: Yes, order 1.

Consider example

$\frac{dy}{dx} = -y$ with $y = y_0$ at $x = 0$

Q: Exact soln is?

A: $y = y_0 e^{-x}$

Euler gives

$f_0 = -y_0$

$y_1 = y(x + \Delta x) = y_0 - y_0 \Delta x = y_0 (1 - \Delta x)$

Repeat: at n'th step

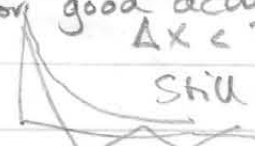
$y_n = y(x + n\Delta x) = y_0 (1 - \Delta x)^n$



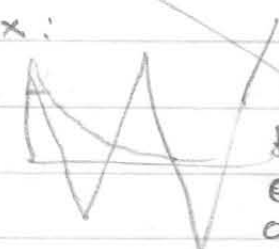
Behavior depends on step size Δx:
 clearly want Δx < 1 for good accuracy



OK soln



Still OK



Δx > 2
 Blowup
 exply,
 oscillating
 Ouch!

Q: Why would you be so stupid as to choose Δx > 2?

A: You've probably got a more complicated set of coupled eqs such as

$\frac{d\vec{y}}{dx} = \underset{\substack{\uparrow \\ \text{matrix}}}{A} \vec{y}$ i.e. $\begin{pmatrix} y_1' \\ \vdots \\ y_n' \end{pmatrix} = \begin{pmatrix} A_{11} & \dots & A_{1n} \\ \vdots & & \vdots \\ A_{n1} & \dots & A_{nn} \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$

For example, y_i might be:

- pops of energy levels of atom
- nuclear abundances in reaction chain
- ionization species
- molecular abundances in chemical chain

Typically "rate coefficients" A_{ij} depend on ^{eg} temperature and other global quantities, but for present purposes treat them as constant.

Q: ^{How to derive} What is exact soln if A_{ij} are const.?

A: Diagonalize $A = P^{-1} \Lambda P$

Then $X = Y P y$ ↑
diagonal mx, whose elts
are eigenvalues of Λ

$$\frac{dy}{dx} = P^{-1} \Lambda P y$$

Define $Y = P y$ ← eigenvectors of Λ

$$\text{Then } \frac{dY}{dx} = \Lambda Y$$

$$\text{ie } \frac{dY_i}{dx} = \Lambda_i Y_i$$

$$\text{solution } Y_i = Y_i(x_0) e^{\Lambda_i(x-x_0)}$$

Problem with this is that some Λ_i may be tiny compared to others.

Y_i with tiny Λ_i are hardly changing;
 Y_i with large +ve Λ_i are increasing exply,
 Y_i with -ve Λ_i are decaying exply.

Numerical scheme sees that latter Y_i have decayed to almost zero,

so it increases the step size.
Then bang, your tiny decaying y_i
start blowing up exponentially & oscillating.

Stepsize is limited by the fastest rate of
small λ_i represent slowest largest -ve λ_i .
If increase stepsize above this, then
solution will blow up. Trouble with this
is you end up taking vast numbers of
steps in which most y_i hardly change
at all. Not good.

Q: What to do about this?

A: First thing is, diagnose the problem.

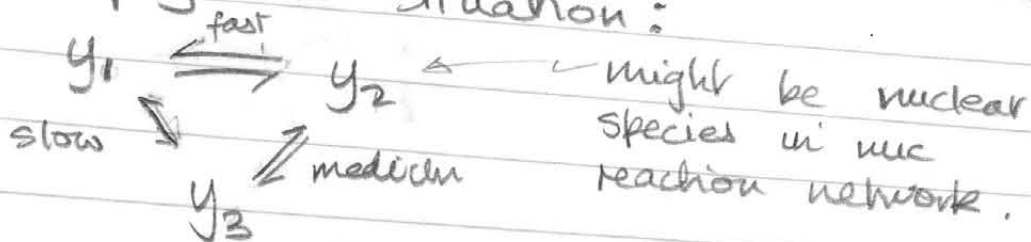
Having diagnosed problem, two
generic options: the "stiff" y_i

- remove rapidly varying y_i from
integration and "fix" them somehow;
- use "stiff" numerical integration method,
will (discussed later).

[Diff eqs like this are called stiff!]

How to apply option (a)?

Typical physical situation:



Fast rates (here $y_0 \rightleftharpoons y_1$) tend to drive y_1 & y_2 into relative "equilibrium".
 So solution to instability is to force set y_1 & y_2 into relative equilibrium.
 Have to implement criterion to do this.

Maybe: \nearrow abundance of species

(1) Is y_1 small, therefore unimportant?

(2) Is stepsize becoming large enough that $y_1 \rightleftharpoons y_2$ is fast

ie $\Delta t \lambda > 1$?

\uparrow appropriate eigenvalue of $\begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$ submatrix.

If so then set y_1 into equilib with y_2 .

Problem of instability affects all "explicit" integration schemes, such as Runge-Kutta.

Explicit = calculate ^{new} value explicitly from earlier values

as opposed to

Implicit = calculate new value from eq involving new value as well as earlier.

↓ Postpone?

Consider, as next most complicated case, 2nd order RK

Recall: Solves $\frac{dy}{dx} = f(x, y)$

from $f_0 = f(x, y)$

$f_1 = f(x + \Delta x, y + f_0 \Delta x)$

$y(x + \Delta x) = y + \frac{1}{2}(f_0 + f_1)\Delta x$

Example (again):

$$\frac{dy}{dx} = -y$$

starting from $y = y_0$
at $x = 0$.

2nd R-K gives

$$f_0 = -y_0$$

$$f_1 = -y_0(1 - \Delta x)$$

so ^{goes into f_1} $y_0 + f_0 \Delta x = y_0(1 - \Delta x)$

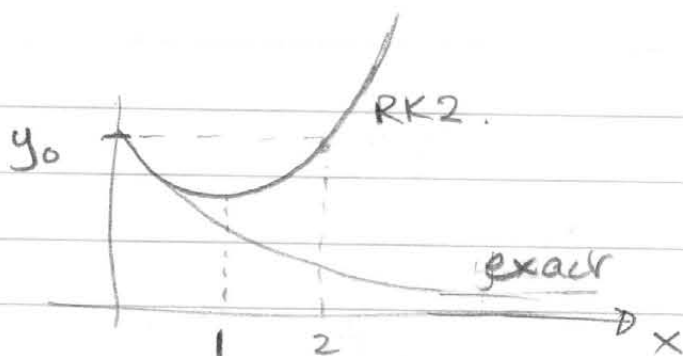
$$y_1 = y(x + \Delta x) = y_0 \left(1 - \Delta x + \frac{\Delta x^2}{2} \right) = y_0 \left[\frac{1}{2} + \frac{1}{2}(1 - \Delta x)^2 \right]$$

whereas true solution is

$$y(x + \Delta x) = y_0 e^{-\Delta x} = y_0 \left(1 - \Delta x + \frac{\Delta x^2}{2!} - \dots \right)$$

At $x = 0$ $y = y_0$

$$y_1 = y_0 \left(1 - \Delta x + \frac{\Delta x^2}{2} \right)$$



Solution decays for $|\Delta x| < 2$,
 but blows up exply + osc for $|\Delta x| > 2$.
 For $|\Delta x| > 2$, blowup is more violent
 than 1st order scheme.

Higher order integrators are more unstable
 than lower order

Stiff integration methods

Are based on implicit schemes

explicit methods

Implicit 1st order RK (= Implicit Euler)

Instead of

explicit: $y_1 = y_0 + f(x_0, y_0) \Delta x$
 solve

implicit: $y_1 = y_0 + f(x_1, y_1) \Delta x$

where you are heading.

Called implicit because solution for y_1 depends on $f(x, y)$ at $y = y_1$.

Could for example try iterative soln.

Clearly more work than explicit, especially if evaluating f is expensive.

Example $y' = -y$

stiff eq.

Instead of

explicit: $y_1 = y_0 - y_0 \Delta x$
 ie $y_1 = y_0(1 - \Delta x)$

solve

implicit: $y_1 = y_0 - y_1 \Delta x$

In this case $y_1 = \frac{y_0}{1 + \Delta x}$

(That was easy, but in general it may not be so easy to solve..)

At n'th step

$$y_n = \frac{y_0}{(1 + \Delta x)^n}$$

This is stable $\forall \Delta x$. Good.

But how does well does implicit method work for growing instead of decaying solutions?

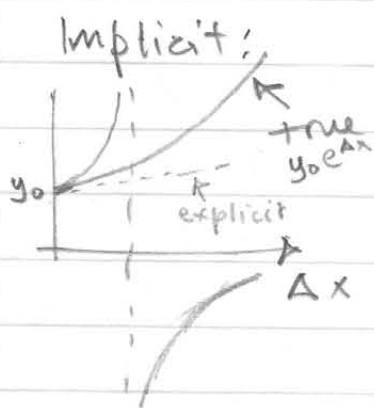
Ex/ $y' = \uparrow y$
 + instead of - .

Here

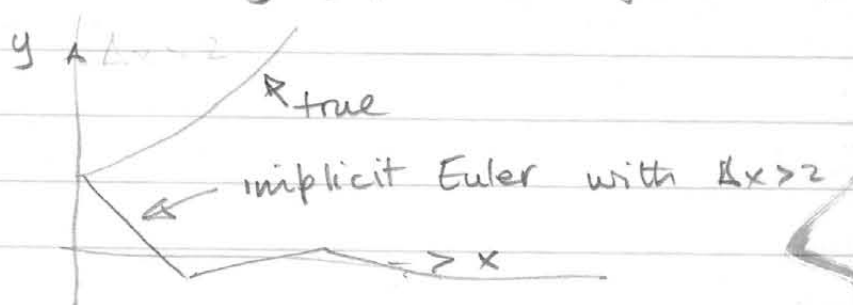
explicit: $y_1 = y_0(1 + \Delta x)$

whereas

implicit: $y_1 = \frac{y_0}{1 - \Delta x}$



Soln grows if $\Delta x < 1$
 grows & osc if $1 < \Delta x < 2$
 decays (!) & osc if $\Delta x > 2$.



Explicit advertized too large stepsize for rapidly decaying modes by blowing up.
 Implicit conceals too large stepsize for rapidly growing modes.

Which poison is worse?

Implicit is safe provided that all modes decaying

Q: Does that ever happen?

A: Yes: approach to thermodynamic equilib.

Before addressing what to do about that instability, let's worry about how to solve implicit eqns.

Linearly implicit Euler

Suppose diffeqs are approx linear:

$$\frac{d\vec{y}}{dx} = A \vec{y}$$

↑
constant mx

ie ~~$\frac{dy_i}{dx} = A_{ij} y_j$~~ (y_i now denotes i^{th} y)
not i^{th} step

explicit Euler;

~~$y_i(x + \Delta x) = y_i(x) + A_{ij} y_j(x) \Delta x$~~
or $\vec{y}(x + \Delta x) = \underbrace{(1 + A \Delta x)}_{\text{matrix}} \vec{y}(x)$

implicit Euler:

$$\vec{y}(x + \Delta x) = \vec{y}(x) + A \vec{y}(x + \Delta x) \Delta x$$

ie $(1 - A \Delta x) \vec{y}(x + \Delta x) = \vec{y}(x)$
ie $\boxed{\vec{y}(x + \Delta x) = \underbrace{(1 - A \Delta x)^{-1}}_{\text{matrix inverse}} \vec{y}(x)}$

So here soln of implicit eqn involves inverting a matrix.

Clearly more expensive numerically than explicit, but not ridiculous by.

In general matrix A is not constant. Linear implicit Euler approximates A by its value at x . Have

$$\frac{dy_i}{dx} \equiv f_i = A_{ij} y_j$$

so $A_{ij} = \partial f_i / \partial y_j$ is "Jacobian" mx

So with this method, derivative routine should also return not only $f_i(x, y)$, but also Jacobian $\text{mx} \frac{\partial f_i}{\partial y_j}(x, y)$.

Obviously ~~this~~ involves more coding overhead than just coding f_i . But sometimes that's OK, because ^{eg} A_{ij} may ^{just} be rate coeffs (in nuclear / atomic / ion / molecular network).

Back to stability issue...

Time symmetric approach

Consider ~~it~~ again

$$y' = Ay \quad (\text{just one } y, \text{ for simplicity})$$

explicit:

$$y_1 = y_0 (1 + A \Delta x)$$

problem: blows up for A -ve and $|A| \Delta x > 2$

implicit:

$$y_1 = \frac{y_0}{1 - A \Delta x}$$

problem: dies for A +ve and $A \Delta x > 2$.

Problem with implicit is just ~~own~~ time-reversed version of explicit.

\Rightarrow idea: what about time symmetric method?

$$y_1 = \left(\frac{1 + A \Delta x / 2}{1 - A \Delta x / 2} \right) y_0$$

Q: is this time symmetric?

A: Yes. Eqn is

$$y_0 = \left(\frac{1 + A \Delta x / 2}{1 - A \Delta x / 2} \right) y_1$$

which is same with $y_0 \leftrightarrow y_1$
and $\Delta x \rightarrow -\Delta x$.

Notice that method is actually 2nd order:

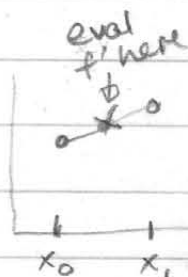
$$\left(\frac{1 + A \Delta x / 2}{1 - A \Delta x / 2} \right) = 1 + A \Delta x + \frac{(A \Delta x)^2}{2!} + O(\Delta x^3)$$

$$= e^{A \Delta x} \text{ to 2nd order.}$$

Simplest time-sym method:

Implicit midpoint Euler

$$y_1 = y_0 + f\left(\frac{x_0 + x_1}{2}, \frac{y_0 + y_1}{2}\right) \Delta x$$



Ex/ $y' = Ay$

$$y_1 = y_0 + A \left(\frac{y_0 + y_1}{2} \right) \Delta x$$

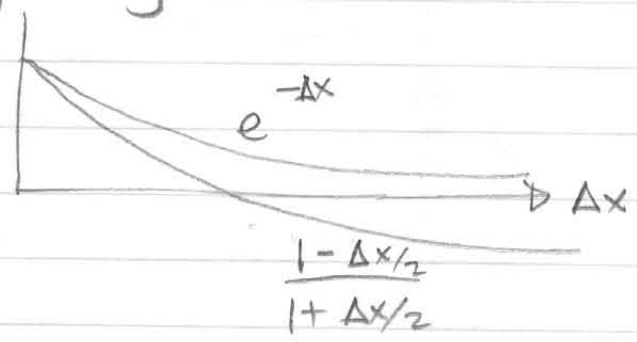
$$\text{ie } \left(1 - \frac{A \Delta x}{2} \right) y_1 = y_0 \left(1 + \frac{A \Delta x}{2} \right)$$

$$\text{ie } y_1 = \frac{\left(1 + \frac{A \Delta x}{2} \right)}{\left(1 - \frac{A \Delta x}{2} \right)} y_0 \text{ as desired.}$$

At nth step

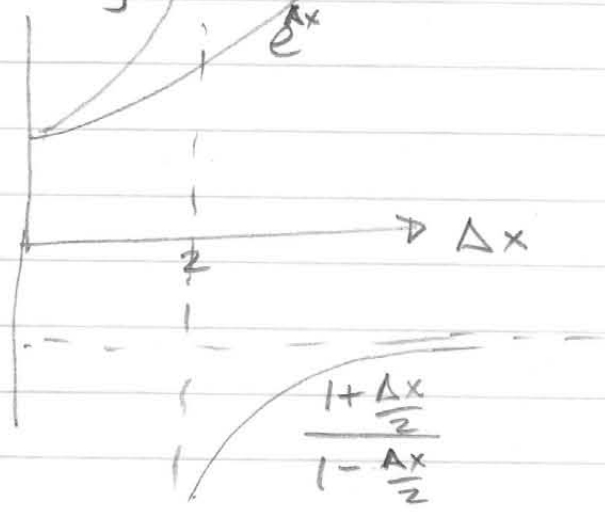
$$y_n = \left(\frac{1 + A \Delta x / 2}{1 - A \Delta x / 2} \right)^n y_0$$

Stability?
 $y' = -y$



Stable $\forall \Delta x$.
 As $\Delta x \rightarrow \infty$
 $y_n \rightarrow (-1)^n y_0$
 oscillates
 with neutral stability.
 OK.

$y' = y$

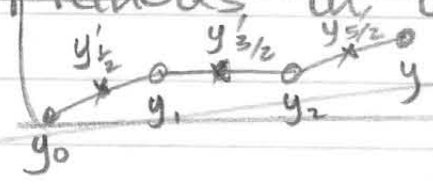


For $\Delta x > 2$,
 oscillates,
 but grows $\rightarrow \infty$
 At least it
 is not hiding
 growth.

Conclusion: time-symmetric approach
 is safer, and generally preferred.

Implicit midpoint Euler is also nice
 in that it is 2nd order with
 even though f is evaluated at 1 pt (midpoint),
 the

Methods in which y' is evaluated only
 at midpoints of y
 are called "leap-frog".



Linearly implicit midpoint Euler
Approximates diffeqs as linear

$$\frac{d\vec{y}}{dx} = A \vec{y}$$

↑
matrix.

Implicit mid Euler gives

$$y_1 = \underbrace{\left(1 - A \frac{\Delta x}{2}\right)^{-1}}_{\text{matrix}} \underbrace{\left(1 + A \frac{\Delta x}{2}\right)}_{\text{matrix}} y_0$$

but no need for mx op here.

But $\left(1 + A \frac{\Delta x}{2}\right) y_0 \approx y_0 + A y_0 \frac{\Delta x}{2}$
(not $y_0 + f_0 \frac{\Delta x}{2}$)
= $y_0 + f_0 \frac{\Delta x}{2}$

So $y_1 = \left(1 - A \frac{\Delta x}{2}\right)^{-1} \left(y_0 + f_0 \frac{\Delta x}{2}\right)$

$$f_0 = f(x_0, y_0)$$
$$y_1 = \left(1 - \underbrace{A \frac{\Delta x}{2}}_{\text{Jacobian } \frac{\partial f}{\partial y} \text{ at } x_0}\right)^{-1} \left(y_0 + f_0 \frac{\Delta x}{2}\right)$$

Jacobian $\frac{\partial f}{\partial y}$ at x_0

Higher-order implicit methods

Can be quite complicated

eg Kaps - Rentrop (eg *Times* 1999)
Bader - Deufhard (ApJS, 124, 241)

Stepsize control for integrator

Integrator should

- Compute to specified accuracy (relative or abs)
- Adjust stepsize to achieve required accuracy.

Common approach is to compare results at stepsize Δx and $2\Delta x$, use that as estimate of error.

Ex / 4th order integrator (eg RK).

One step of length $2\Delta x$:

$$y_1(x + 2\Delta x) = y_{true} + \overbrace{(2\Delta x)^5}^{\text{error in } y} a + O(\Delta x^6)$$

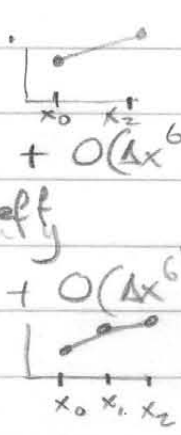
Two steps of length Δx :

$$y_2(x + 2\Delta x) = y_{true} + 2 \overbrace{(\Delta x)^5}^{\text{some coeff error}} a + O(\Delta x^6)$$

So estimate of error is

$$\Delta y_2 \approx \frac{1}{2} (2(\Delta x)^5 a) = \Delta x^5 \sigma$$

$$\text{ie } \sigma = \frac{y_2 - y_1}{2^4 - 1} \Delta x = \frac{y_2 - y_1}{15} \Delta x^5$$



If you want error $\sigma_{desired}$, set $\frac{\Delta y_2}{\Delta x} = \left(\frac{\Delta x}{\Delta x_{desired}} \right)^5$

$$\text{or } \Delta x_{desired} = \Delta x \left(\frac{\sigma}{\Delta y_2} \right)^{1/5}$$

More typically, change Δx by factor of 2, doubling or halving Δx as nec.

Thus if $\Delta y_2 > \frac{1}{2} \sigma_{desired}$, then halve Δx

$$\Delta y_2 < \frac{1}{25} \sigma_{desired}, \text{ then double } \Delta x$$

But in practice you may want different errors σ_i for different y_i .
 Typically, you may want
 large $|y_i|$: $\sigma_i \approx \epsilon |y_i|$ = constant relative err
 small $|y_i|$: $\sigma_i \approx \epsilon$ = constant error.

Choose $\sigma_{desired} = \text{smallest of individual } \sigma_{i,desired}$.

Improved estimate

As by-product of need to estimate error, can improve estimate of y_{true} by combining y_1 & y_2 :

$$\frac{2^4 y_2 + y_1}{2^4 - 1} = y_{true} + O(\Delta x^6)$$

$$= \frac{16}{15} y_2 - \frac{1}{15} y_1$$

Gained an extra power of Δx in accuracy!

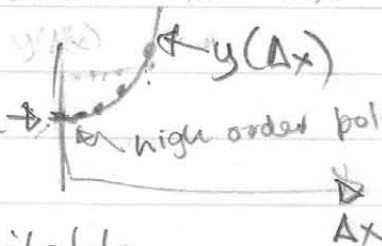
Notion of using multiple stepsizes

(a) to estimate error, (b) to reduce error leads to: different stepsizes

Richardson extrapolation

Idea: use several stepsizes Δx over same interval, then extrapolate to $\Delta x \rightarrow 0$.
 "sufficiently analytic"

For some systems, y_{true} of ODEs, this approach is best available,



But warned: polyn extrap can fail,

and rat fn extrap can fail.
 and fn to be fitted should be ^{so} order polyn.