

## Other Krylov subspace methods

Idea: find "best" updates inside Krylov subspaces for  $A$  (or sometimes  $A^T$ ). Some methods work even with unsymmetric problems. Here are some examples:

### GMRES (Generalized minimal residual)

→ works for general systems, residual guaranteed to decrease but iterations become more costly (in computations / storage)

Restarted GMRES: run  $k$  steps of GMRES, throw away "memory" and start  $k$  steps of GMRES again from where we left over.

CGN: Conjugate Gradient for normal eq: specially designed to solve  $ATAx = ATb$ .

BCG: Bi-Conjugate gradient (for general problems)  
(no optimality and unstable)

Bi-CGSTAB: stabilized version

QMR: quasi-minimal residual (another stabilization of BCG)

etc..

- Krylov subspace methods are also used for finding eigenvalues (we shall discuss this later)
- Matlab has a good collection of methods for solving linear systems:  
    help spfun

## Chap 5 Initial value problems for Ordinary Diff Eq (ODE)

- Objective develop methods to solve numerically problems of the kind

$$(1) \quad \left\{ \begin{array}{l} \frac{dy}{dt} = f(t, y) \quad \text{for } t \in [a, b] \\ y(a) = \alpha \end{array} \right. \quad \left\{ \begin{array}{l} \text{often an exact solution to (1)} \\ \text{is too complicated or even} \\ \text{cannot be found.} \end{array} \right.$$

Extension to systems:

$$\left\{ \begin{array}{l} \frac{dy_1}{dt} = f_1(t, y_1, y_2, \dots, y_n) \\ \frac{dy_2}{dt} = f_2(t, y_1, y_2, \dots, y_n) \\ \vdots \\ \frac{dy_n}{dt} = f_n(t, y_1, y_2, \dots, y_n) \end{array} \right. \quad \begin{aligned} & \text{for } t \in [a, b], \text{ s.t.} \\ & y_i(a) = \alpha_i, i = 1 \dots n \end{aligned}$$

$$\left( \frac{dy}{dt} = f(t, y) \right)$$

And to  $n$ -th order IVP:

$$y^{(n)} = f(t, y, y', \dots, y^{(n-1)}) \quad \text{s.t.} \quad \left\{ \begin{array}{l} y(a) = \alpha_1 \\ y'(a) = \alpha_2 \\ \vdots \\ y^{(n-1)}(a) = \alpha_n \end{array} \right.$$

⚠ this is different from notion of stability for numerical methods

Fundamental questions to answer about the IVP (1)

- existence: does (1) admit a sol?
- uniqueness: is the sol to (1) unique?
- "stability": Do small changes in the statement of the problem introduce small changes in the solution?

A problem satisfying all 3 properties is said to be well posed  
(in the sense of Hadamard)

The IVP (1) is well posed under mild assumptions on  $f$ . We need some terminology (maybe you've seen this many times before!) (2) 35

### Def (Lipschitz condition)

A function  $f(t, y)$  satisfies a Lipschitz condition in the variable  $y$  on some  $D \subset \mathbb{R}^2$  iff  $\exists L > 0$

$$|f(t, y_1) - f(t, y_2)| \leq L |y_1 - y_2| \quad \forall (t, y_1) \in D \\ (t, y_2) \in D.$$

$L$  = Lipschitz constant for  $f$ .

(means  $f(t, \cdot)$  is Lipschitz continuous for all  $t$ ).

### Example:

$$D = \{(t, y) \mid 0 \leq t \leq 1, -1 \leq y \leq 1\}$$

$$f(t, y) = t|y|$$

$$|f(t, y_1) - f(t, y_2)| = |t(y_1 - t|y_2|)| = |t||y_1 - y_2| \leq |y_1 - y_2|$$

$$L = 1$$

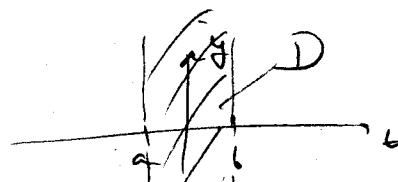
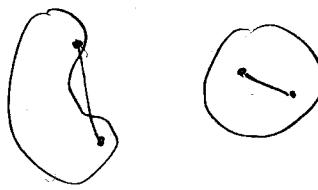
### Def (Convex set) A set $D$ is a convex set iff:

$$x, y \in D \Rightarrow \lambda x + (1-\lambda)y \in D \quad \forall \lambda \in [0, 1]$$

in our case:

$$(t_1, y_1) \in D, (t_2, y_2) \in D \rightarrow (\lambda t_1 + (1-\lambda)t_2, \lambda y_1 + (1-\lambda)y_2)$$

$$\forall \lambda \in [0, 1]$$



For IVP we usually have:

$$D = \{f(t, y) \mid a \leq t \leq b, y \in \mathbb{R}\}$$

which is convex

Here is a sufficient (but not necessary) condition for Lipschitz cond.: (3)  
36

Theorem: Let  $f(t, y)$  be defined on some convex set  $D \subset \mathbb{R}^2$ .

If  $\exists L > 0$  s.t.

$$\left| \frac{\partial f}{\partial y}(t, y) \right| \leq L \quad \forall (t, y) \in D$$

then  $f$  satisfies a Lipschitz condition on  $D$  on  $y$  with Lip. const.  $L$ .

Note:  $f(t, y) = t|y|$  satisfies the Lipschitz condition and yet

$\frac{\partial f}{\partial y}$  does not exist at  $y=0$ .

Existence and uniqueness for IVP (C) are taken care by:

Theorem : Let  $D = \{(t, y) \mid a \leq t \leq b \text{ & } y \in \mathbb{R}\}$  and that

- $f(t, y)$  is continuous on  $D$
  - $f$  satisfies Lip. cond on  $D$  in Variable  $y$ , then the IVP(1)  
admits a unique solution.

"Stability": can be formulated as follows:

$$\text{Use } y(t) \text{ solve } \begin{cases} \frac{dy}{dt} = f(t, y), & a \leq t \leq b \\ y(a) = \alpha \end{cases}$$

$\exists \epsilon_0 > 0, k > 0 \quad \forall \epsilon \text{ s.t. } \epsilon_0 > \epsilon > 0$

If  $\delta(t)$  continuous s.t.  $|\delta(t)| \leq \varepsilon$

$$\forall \delta_0 \in \mathbb{R} \quad \text{s.t. } |\delta_0| \leq \varepsilon$$

$$\begin{cases} \frac{d\beta}{dt} = f(t, \beta) + \delta(t), & a \leq t \leq b, \\ \beta(a) = \alpha + \delta_0. \end{cases} \quad (\text{perturbed problem})$$

has a unique solution  $z(t)$  and

$$|z(t) - y(t)| < k \epsilon$$

(Essentially the mapping data  $(\alpha, f(t, y))$  to solution  $y$  is continuous).

This "stability" property is crucial to trust solutions given by a numerical method (sources of error can be from the method itself or from numerical roundoff).

~ All numerical methods assume IVP is well-posed.

Example:  $D = \{(t, y) \mid t \in [0, 2], y \in \mathbb{R}\}$

$$\begin{cases} \frac{dy}{dt} = y - t^2 + 1, & 0 \leq t \leq 2 \\ y(0) = \frac{1}{2} \end{cases} \quad (\text{IVP})$$

$$\left| \frac{\partial f}{\partial y} \right| = |1| = 1 \Rightarrow f(t, y) \text{ satisfies Lip. cond on } D \text{ with variable } y.$$

&  $f$  constant  $\Rightarrow$  problem is stable to perturbations in init. data.

We can verify this directly:

$$\begin{cases} \frac{dz}{dt} = z - t^2 + 1 + \delta & \text{(perturbed problem)} \\ z(0) = \frac{1}{2} + \delta_0 \end{cases}$$

$$\text{IVP has sol: } \begin{aligned} y(t) &= \frac{-1}{2} e^t + (t+1)^2 \\ z(t) &= \left(\frac{-1}{2} + \delta_0\right) e^t + (t+1)^2 + (e^t - 1)\delta \end{aligned}$$

$$|z - y| = |((\delta_0 + \delta)e^t - \delta)| \leq |\delta_0 + \delta| e^2 + |\delta| \leq 2e^2 \epsilon + \epsilon = (2e^2 + 1)\epsilon$$

## § 5.2 Euler's method

A simple numerical method to solve IVP:

38

$$\begin{cases} \frac{dy}{dt} = f(t, y), & a \leq t \leq b \\ y(a) = \alpha \end{cases} \quad (\text{assuming well-posedness})$$

Euler's method gives approx only on some mesh points:

$$t_j = a + h j \quad \text{where } h = \frac{b-a}{N} = \text{step size}$$

$$t_0 = a, t_1, t_2, \dots, t_{N-1}, t_N = b$$

( $N+1$  equally spaced points)  
(can be relaxed)

Idea for Euler's method: Taylor's theorem:

$$y(t_{i+1}) = y(t_i) + y'(t_i) \underbrace{(t_{i+1} - t_i)}_h + \frac{1}{2} y''(\xi_i) \underbrace{(t_{i+1} - t_i)^2}_h^2$$

for some  $\xi_i \in (t_i, t_{i+1})$ .

$$= y(t_i) + h y'(t_i) + \frac{h^2}{2} y''(\xi_i)$$

$$= y(t_i) + h f(t_i, y(t_i)) + \underbrace{h^2 y''(\xi_i)}_{\substack{\text{neglect for} \\ \text{Euler's method}}} \quad \begin{array}{l} \text{($y$ satisfies} \\ \text{DE in IVP}) \end{array}$$

$$y_0 = \alpha$$

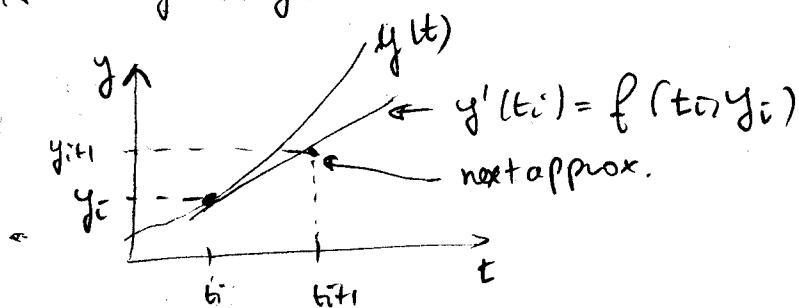
for  $i = 0 \dots N-1$

known

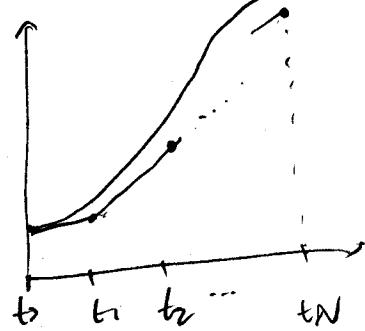
$$y_{i+1} = y_i + h f(t_i, y_i)$$

Geometrical interpretation

Assume  $y_i = y(t_i)$



systematic error  
is introduced at  
many times: ↓  
every step



## Several types of errors in numerical methods for DE:

(6)

39

- Local truncation error: error made in one step-

For example in Euler's method:  $O(h)$  since:  
 book definition w.r.t  $t$ )

$$y(t_i + h) = y(t_i) + h f(t_i, y(t_i)) + O(h^2)$$

- Local roundoff error: precision used for computation ( $10^{-16}$  double  
 $10^{-8}$  single)

- Global truncation error: accumulation of all the local truncation errors. If local truncation error was  $O(h^{n+1})$  then the global truncation error must be  $O(h^n)$  because the number of steps  $\Theta(h)$ .

- Global roundoff error: accumulation of local roundoff errors of previous steps.

- Total error = global truncation error + global roundoff error

If the global truncation error is  $O(h^m)$  then the method is of order m

Example: Euler's method is of order 1, and it is relatively simple to derive more precise bounds on the global truncation error.

Then let  $f$  be continuous and satisfy Lipschitz condition with L-constant L on  $D = \{(t, y) \mid t \in [a, b], y \in \mathbb{R}\}$  and that  $|y''(t)| \leq M \quad \forall t \in [a, b]$  for some  $M > 0$ .

Let  $y(t)$  be the sol to IVP  $\begin{cases} \frac{dy}{dt} = f(t, y) \\ y(a) = \alpha \end{cases}$

and  $y_0, y_1, \dots, y_n$  the approximations given by Euler's method

then :

$$|y(t_i) - y_i| \leq \frac{hM}{2L} [e^{L(t_i-a)} - 1], \quad i=0, 1, \dots, N.$$

proof (sketch)

$$y(t_{i+1}) = y(t_i) + h f(t_i, y(t_i)) + \frac{h^2}{2} y''(\xi_i)$$

$$y_{i+1} = y_i + h f(t_i, y_i)$$

$$\Rightarrow |y(t_{i+1}) - y_{i+1}| \leq |y(t_i) - y_i| + h |f(t_i, y(t_i)) - f(t_i, y_i)| + \frac{h^2}{2} |y''(\xi_i)| \\ \leq (1+hL) |y(t_i) - y_i| + \frac{h^2 M}{2}$$

using Lemma 5.8 in book: ( $\sim$  geometric series + bounding exp)

$$|y_{i+1} - w_{i+1}| \leq e^{(i+1)hL} \left( \underbrace{|y_0 - w_0|}_{=0} + \frac{h^2 M}{2hL} \right) - \frac{h^2 M}{2hL} \\ \leq \frac{h^2 M}{2hL} (e^{(i+1)hL} - 1)$$

$$t_{i+1} - t_0 = t_{i+2} - a$$

Lemma 5.8  $\{a_i\}_{i=0}^k$  is a seq with  $a_0 \geq -\frac{t}{s}$  and

$$a_{i+1} \leq a_i(1+s) + t \quad \text{for } i=0, \dots, k-1$$

$$\Rightarrow a_{i+1} \leq e^{(i+1)s} (a_0 + \frac{t}{s}) - \frac{t}{s}$$

$$\underline{\text{proof}}: a_{i+1} \leq a_i(1+s) + t \leq ((1+s)a_{i-1} + t)(1+s) + t$$

$$\leq ((1+s)((1+s)a_{i-2} + t) + t)(1+s) + t$$

$$\vdots \\ \leq (1+s)^{i+1} a_0 + t \sum_{j=0}^i (1+s)^j$$

$$\Rightarrow a_{i+1} \leq ((1+s)^{i+1} a_0 + t) \left[ \frac{1 - (1+s)^{i+1}}{1 - (1+s)} \right] \quad (3)$$

41

$$= (1+s)^{i+1} \left( a_0 + \frac{t}{s} \right) - \frac{t}{s}$$

$$\leq e^{(i+1)s} \left( a_0 + \frac{t}{s} \right) - \frac{t}{s}$$

$$\text{since } (1+x)^x = e^{x \ln(1+x)} \leq e^x$$

<sup>P</sup> can be shown using Taylor's theorem.

To Euler's method it's even possible to incorporate the round off errors in the analysis:

instead of having:

$$y_0 = a$$

for  $i = 0..N-1$ ,

$$| y_{i+1} = y_i + h f(t_i, y_i) |$$

we commit a mistake at each step:

$$y_0 = a + \delta_0$$

for  $i = 0..N-1$

$$| y_{i+1} = y_i + h f(t_i, y_i) + \delta_{i+1} |$$

If  $|\delta_i| < \delta$  it is possible to show:

$$| y(t_i) - y_0 | \leq \frac{1}{L} \left( \frac{hM}{2} + \frac{\delta}{h} \right) [e^{L(t_i-a)} - 1] + |\delta| e^{L(t_i-a)}$$

performance degrades as  $h \rightarrow 0$  !!

(similar to problems with numerical differentiation)

The  $h$  giving smallest error is

$$h = \sqrt{\frac{2\delta}{M}} \quad (\text{simple calculation})$$

### §3 Higher order Taylor methods

42

Same derivation as Euler method but carry Taylor series further:

$$y(t_{i+1}) = y(t_i) + h y'(t_i) + \frac{h^2}{2} y''(t_i) + \dots + \frac{h^n}{n!} y^{(n)}(t_i) + \frac{h^{n+1}}{(n+1)!} y^{(n+1)}(\xi_i)$$

for some  $\xi_i \in (t_i, t_{i+1})$ .

$$y'(t) = f(t, y(t))$$

$$y''(t) = \frac{d}{dt}[f(t, y(t))]$$

$$y^{(k)}(t) = \frac{d^{k-1}}{dt^{k-1}} [f(t, y(t))]$$

### Taylor method of order n

$$y_0 = \alpha$$

for  $i=0 \dots N-1$ ,

$$y_{i+1} = y_i + h f(t, y(t)) + \frac{h^2}{2} \frac{d}{dt}[f(t, y(t))] + \dots + \frac{h^n}{n!} \frac{d^{n-1}}{dt^{n-1}} [f(t, y(t))]$$

Local truncation error is  $O(h^{n+1})$

In general evaluating  $\frac{d^k}{dt^k} [f(t, y(t))]$  can be quite involved because of the repeated application of the chain rule.

For example:

$$\begin{cases} y' = \cos t - \sin y + t^2 \\ y(-1) = 3 \end{cases} = f(t, y)$$

$$y'' = \frac{d}{dt}[f(t, y(t))] = -\sin t - y' \cos y + 2t$$

$$y''' = \frac{d^2}{dt^2}[f(t, y(t))] = -\cos t - y'' \cos y + (y')^2 \sin y + 2$$

$$y^{(4)} = \frac{d^3}{dt^3}[f(t, y(t))] = \sin t - y^{(3)} \cos y + 3y'y'' \sin y + (y')^3 \cos y$$

etc...

Note: • Assumes  $f$  is smooth ( $n-1$  times differentiable iff method of order  $n$  is desired)

- The higher the accuracy the more complicated the steps become
- differentiation can be carried out symbolically if possible!

#### § 5.4 Runge-Kutta methods

Taylor series method for  $\begin{cases} \frac{dy}{dt} = f(t, y), t \in [a, b] \\ y(a) = A \end{cases}$

require us to compute formulas for  $y'' = \frac{d^2f(t, y)}{dt^2}$   
 $y^{(3)} = \frac{d^3f(t, y)}{dt^3}$   
etc...

which can be quite involved.

Runge-Kutta methods avoid this difficulty by carefully chosen combinations of values of  $f(t, y)$ .

#### Second order Runge-Kutta methods:

Start with Taylor series:

$$y(t+h) = y(t) + h y'(t) + \frac{h^2}{2} y''(t) + \frac{h^3}{3!} y^{(3)}(t) + \dots$$

From the DE we get:

$$y'(t) = f(t, y) = f$$

$$y''(t) = \frac{d}{dt} f(t, y) = \frac{\partial f}{\partial t}(t, y) + y' \frac{\partial f}{\partial y}(t, y) = f_t + f f_y$$

$$y'''(t) = \frac{d^2}{dt^2} f(t, y) = f_{tt} + f f_{ty} + (f_t + f f_y) f_y + f (f_{ty} + f f_{yy})$$

etc...

$$\Rightarrow y(t+h) = y(t) + hf + \frac{h^2}{2}(f_t + f f_y) + O(h^3)$$

The idea is to eliminate the partial derivative of  $f$  by using the first few terms of the two variable Taylor series for  $f(t, y)$ :

$$\begin{aligned} f(t+h, y+hf) &= f + Df \cdot \begin{pmatrix} h \\ hf \end{pmatrix} + O(h^2) \\ &= f + h f_t + h f f_y + O(h^2) \end{aligned}$$

Rewriting the Taylor series of  $y$ :

$$\begin{aligned} (*) \quad y(t+h) &= y(t) + \frac{1}{2} h f + \frac{1}{2} h [f + h f_t + h f f_y] + O(h^3) \\ &\quad \text{first 2 terms of 2 variable Taylor series of } f. \\ &= y(t) + \frac{1}{2} h f + \frac{1}{2} h [f(t+h, y+hf) + O(h^2)] \\ &\quad + O(h^3). \end{aligned}$$

Thus it's possible to construct an update which has the same  $O(h^3)$  local truncation error as 2nd-order Taylor method:

### Modified Euler method

$$y_0 = A$$

for  $i = 0, 1, \dots, N-1$

$$F_1 = h f(t_i, y_i)$$

$$F_2 = h f(t_i + h, y_i + F_1)$$

$$y_{i+1} = y_i + \frac{1}{2} F_1 + \frac{1}{2} F_2$$

The general form for Second order Runge Kutta updates is:

$$y(t+h) = y + w_1 h f + w_2 h f(t+\alpha h, y+\beta h f)$$

where  $w_1, w_2, \alpha, \beta$  are parameters that we can adjust.

Using two variable Taylor expansion:

$$y(t+h) = y + w_1 h f + w_2 h [f + \alpha h f_t + \beta h f f_y]$$

Matching comparable terms with (\*) we get:

(12)

45

$$\left\{ \begin{array}{l} w_1 + w_2 = 1 \\ w_2 \alpha = \frac{1}{2} \\ w_2 \beta = \frac{1}{2} \end{array} \right. \rightarrow \text{so there is a family of RK-} \\ \text{order 2 methods}$$

Modified Euler:  $w_1 = w_2 = \frac{1}{2}, \alpha = \beta = 1$

Midpoint method:  $w_1 = 0, w_2 = 1, \alpha = \beta = \frac{1}{2}$

$$y_0 = A$$

for  $i = 0, 1, \dots, N-1$

(also  $O(h^3)$  LTE)

$$F_1 = h f(t_i, y_i)$$

$$F_2 = h f(t_i + \frac{h}{2}, y_i + \frac{1}{2} F_1)$$

$$y_{i+1} = y_i + F_2$$

Huen's method:  $w_1 = \frac{1}{4}, w_2 = \frac{3}{4}, \alpha = \beta = \frac{2}{3}$

$$y_0 = A$$

for  $i = 0, 1, \dots, N-1$

$$F_1 = h f(t_i, y_i)$$

$$F_2 = h f(t_i + \frac{2}{3}h, y_i + \frac{2}{3} F_1)$$

$$y_{i+1} = y_i + \frac{1}{4} F_1 + \frac{3}{4} F_2$$

Runge Kutta methods of order 3 are obtained

by matching terms between Taylor expansion of order 3 of  $y(t)$  46  
and  $f(t + \alpha h, y + w_1 f(t + \beta h, y + w_2 f(t, y)))$

The derivation is quite tedious, but here is one possible RK order 3 method:

$$y_0 = A$$

for  $i = 0, \dots, N-1$

$$F_1 = h f(t_i, y_i)$$

$$F_2 = h f\left(t_i + \frac{1}{2}h, y_i + \frac{1}{2}F_1\right)$$

$$F_3 = h f\left(t_i + \frac{3}{4}h, y_i + \frac{3}{4}F_2\right)$$

$$y_{i+1} = y_i + \frac{1}{6}(2F_1 + 3F_2 + 4F_3)$$

$$\begin{pmatrix} \alpha = w_1 = \frac{3}{4} \\ \beta = w_2 = \frac{1}{2} \end{pmatrix}$$

However it is not commonly used in practice.

The most popular RK method is that of order 4: again the derivation is tedious but the implementation straightforward.

Runge Kutta method of order 4 (LTE  $\mathcal{O}(h^5)$ )

$$y_0 = A$$

for  $i = 0, \dots, N-1$

$$F_1 = h f(t_i, y_i)$$

$$F_2 = h f\left(t_i + \frac{1}{2}h, y_i + \frac{1}{2}F_1\right)$$

$$F_3 = h f\left(t_i + \frac{1}{2}h, y_i + \frac{1}{2}F_2\right)$$

$$F_4 = h f\left(t_i + h, y_i + F_3\right)$$

$$y_{i+1} = y_i + \frac{1}{6}(F_1 + 2F_2 + 2F_3 + F_4)$$

# of function eval.      1 2 3 4 5 6 7 8 9 10 11 12 ...

max order of RK method      1 2 3 4 4 5 6 6 7 7 8 9 ...

So # of function eval increases more rapidly than the max order of RK method.  $\Rightarrow$  higher order RK methods are less attractive than the classical RK4. (it makes more sense to use smaller time steps for RK4 than using a higher order method with bigger steps).

### § 5.5 Adaptive Runge Kutta Fehlberg method

Idea: Estimate local truncation error and adjust the step length accordingly.

#### Local truncation error estimation

Assume we are given two update formulas for solving

$$\begin{cases} y' = f(t, y) & , t \in [a, b] \\ y(a) = \alpha \end{cases}$$

with local truncation errors differing by  $1$ :

$$\textcircled{1} \quad y(t_{i+1}) = y(t_i) + h \phi(t_i, y(t_i), h) + O(h^{n+1}) \quad (\text{order } n)$$

$$y_0 = \alpha$$

for  $i = 0, \dots, N-1$

$$\mid y_{i+1} = y_i + h \phi(t_i, y_i, h)$$

$$\textcircled{2} \quad y(t_{i+1}) = y(t_i) + h \tilde{\phi}(t_i, y(t_i), h) + O(h^{n+2}) \quad (\text{order } n+1)$$

$$\tilde{y}_0 = \alpha$$

for  $i = 0 \dots N-1$

$$\mid \tilde{y}_{i+1} = \tilde{y}_i + h \tilde{\phi}(t_i, y_i, h)$$