

# A Probabilistic Model for the Numerical Solution of Initial Value Problems

by Schober, Särkkä and Hennig

---

Jon Cockayne

July 26, 2017

In Schober et al. [2016]...

- Schober et al. [2014] showed that for certain choices of GP prior, sequential conditioning led to a Runge-Kutta mean.

---

<sup>1</sup>Skeel [1979]

In Schober et al. [2016]...

- Schober et al. [2014] showed that for certain choices of GP prior, sequential conditioning led to a Runge-Kutta mean.
- This could only be applied over a **single time step**.

---

<sup>1</sup>Skeel [1979]

In Schober et al. [2016]...

- Schober et al. [2014] showed that for certain choices of GP prior, sequential conditioning led to a Runge-Kutta mean.
- This could only be applied over a **single time step**.
- In this paper it is shown that for those priors, the predictive posterior corresponds to a **linear multistep method** for solving the ODE in **Nordsieck form**<sup>1</sup>.

---

<sup>1</sup>Skeel [1979]

In Schober et al. [2016]...

- Schober et al. [2014] showed that for certain choices of GP prior, sequential conditioning led to a Runge-Kutta mean.
- This could only be applied over a **single time step**.
- In this paper it is shown that for those priors, the predictive posterior corresponds to a **linear multistep method** for solving the ODE in **Nordsieck form**<sup>1</sup>.
- This allows us to establish **global error results** in a consistent fashion.

---

<sup>1</sup>Skeel [1979]

## Standard Solvers and Nordsieck Form

---

## Problem Setup

We have an IVP:

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

where

- $t \in \mathbb{T} := [t_0, T] \subset \mathbb{R}$
- $y : \mathbb{T} \rightarrow \mathbb{R}$
- $f : \mathbb{T} \times \mathbb{R} \rightarrow \mathbb{R}$ .
- $y(t_0) = y_0$ .

## Problem Setup

We have an IVP:

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

where

- $t \in \mathbb{T} := [t_0, T] \subset \mathbb{R}$
- $y : \mathbb{T} \rightarrow \mathbb{R}$
- $f : \mathbb{T} \times \mathbb{R} \rightarrow \mathbb{R}$ .
- $y(t_0) = y_0$ .

Produce a solution on a **grid**  $\{t_n\}$ ,  $n = 0, 1, 2, \dots, N$ . Let  $y_n$  be the approximation to  $y$  at time  $t_n$ , and  $h_n = t_n - t_{n-1}$ . Let  $z_n = f(t_n, y_n)$ .



Only consider **explicit** methods.

Only consider **explicit** methods.

The most basic: **Explicit Euler**

$$y_n = y_{n-1} + hz_{n-1}$$

Problems: global error  $O(h)$ , unstable.

We would like methods of order  $q$  - that is, with global error  $O(h^q)$

# Linear Multistep Methods

(Explicit) **Linear Multistep Methods** approximate the solution as a linear combination of previous evaluations of  $f$  and estimates of  $y$

$$\sum_{i=0}^k \alpha_i y_{n-i} = h \sum_{i=1}^k \beta_i z_{n-i} \quad (1)$$

(Explicit) **Linear Multistep Methods** approximate the solution as a linear combination of previous evaluations of  $f$  and estimates of  $y$

$$\sum_{i=0}^k \alpha_i y_{n-i} = h \sum_{i=1}^k \beta_i z_{n-i} \quad (1)$$

Two predominant methods:

- Adams-Bashforth
- Adams-Moulton

[See also: Teymur et al. [2016]]

# Runge Kutta Methods

Runge Kutta Methods evaluate  $f$  at multiple locations **between** time-steps:

$$k_i = f(t_n + c_i h, y + h \sum_{j < i} a_{ij} k_j)$$

$$y_{n+1} = y_n + h \sum_{i=1}^s b_i k_i$$

0			
$c_2$	$a_{21}$		
$c_3$	$a_{31}$	$a_{32}$	
	$b_1$	$b_2$	$b_3$

# Runge Kutta Methods

Butcher Tableau for  $q = s = 2$ :

$$k_1 = f(t_n, y_n)$$

$$k_2 = f(t_n + \alpha h, y_n + \alpha h k_1)$$

$$y_{n+1} = y_n + h \left(1 - \frac{1}{2\alpha}\right) k_1 + \frac{h}{2\alpha} k_2$$

Midpoint method:  $\alpha = \frac{1}{2}$

Heun's method:  $\alpha = 1$

0		
$\alpha$	$\alpha$	
	$(1 - \frac{1}{2\alpha})$	$\frac{1}{2\alpha}$

All Linear Multistep Methods<sup>2</sup> can be written in **Nordsieck form**:

$$\mathbf{x}_n = \left( y_n, hy'_n, \dots, \frac{h^q y_n^{(q)}}{q!} \right)$$
$$\mathbf{x}_{n+1} = (\mathbf{I} - \mathbf{L}\mathbf{e}_1^T) \mathbf{P}\mathbf{x}_n + h\mathbf{L}z_n$$

where  $z_n$  solves

$$[\mathbf{P}\mathbf{x}_n]_1 + [\mathbf{L}]_1 z_n = hf(t_n + h, [\mathbf{P}\mathbf{x}_n]_0 + h[\mathbf{L}]_0 z_n)$$

$\mathbf{P}$  the Pascal Triangle matrix.

---

<sup>2</sup>and some Runge-Kutta Methods?

For suitable choices of  $L$ , Nordsieck methods can achieve local truncation error (at least)  $q$ .



## A Probabilistic Model

---

For a probability space  $(\Omega, \mathcal{F}, P)$ ...

- $X_t$  the prior distribution for the solution  $y(t)$ <sup>3</sup>, taking values in  $\mathbb{R}^{q+1}$ .
- $Y$  the *true* solution
- $Y'$  its first derivative.

---

<sup>3</sup>This is somewhat unclear from the text, but I believe it should be considered as the prior

For a probability space  $(\Omega, \mathcal{F}, P)$ ...

- $X_t$  the prior distribution for the solution  $y(t)$ <sup>3</sup>, taking values in  $\mathbb{R}^{q+1}$ .
- $Y$  the *true* solution
- $Y'$  its first derivative.

We also let  $\mathcal{F}_t$  denote the **filtration** generated by  $X_t$ .

---

<sup>3</sup>This is somewhat unclear from the text, but I believe it should be considered as the prior

Observation: certain GPs can be written as a **Linear Time Invariant** (LTI) SDE:

$$dX_t = FX_t dt + LdW_t$$

$F$  the “state feedback matrix” and  $L$  the “diffusion vector”.  $W_t$  a Wiener process,  $dW(t) \sim \mathcal{N}(0, \sigma^2 dt)$ .

Observation: certain GPs can be written as a **Linear Time Invariant** (LTI) SDE:

$$d\mathbf{X}_t = \mathbf{F}\mathbf{X}_t dt + \mathbf{L}dW_t$$

$\mathbf{F}$  the “state feedback matrix” and  $\mathbf{L}$  the “diffusion vector”.  $W_t$  a Wiener process,  $dW(t) \sim \mathcal{N}(0, \sigma^2 dt)$ .

We assume  $Y_t$  and  $Y_t'$  are related to  $\mathbf{X}_t$  by

$$Y_t = \mathbf{H}_0 \mathbf{X}_t$$

$$Y_t' = \mathbf{H}_1 \mathbf{X}_t$$

Given a  $X_{t_*} \sim \mathcal{N}(m_*, C_*)$ , for  $t > t_*$  we have that  $X_t$  is also Gaussian.

## SDE Formulation

Given a  $X_{t_*} \sim \mathcal{N}(m_*, C_*)$ , for  $t > t_*$  we have that  $X_t$  is **also Gaussian**.

Let  $A(h) = \exp(Fh)$ . Then

$$\begin{aligned} m_t &= A(t - t_*)m_* \\ \text{cov}(X_t, X_{t'}) &= A(t' - t_*)C_*A(t - t_*)^\top \\ &\quad + \int_{t_*}^{\min(t, t')} A(t - \tau)L\sigma^2L^\top A(t' - \tau)^\top d\tau \end{aligned}$$

## SDE Formulation

Given a  $\mathbf{X}_{t_*} \sim \mathcal{N}(\mathbf{m}_*, \mathbf{C}_*)$ , for  $t > t_*$  we have that  $\mathbf{X}_t$  is **also Gaussian**.

Let  $\mathbf{A}(h) = \exp(\mathbf{F}h)$ . Then

$$\begin{aligned} \mathbf{m}_t &= \mathbf{A}(t - t_*)\mathbf{m}_* \\ \text{cov}(\mathbf{X}_t, \mathbf{X}_{t'}) &= \mathbf{A}(t' - t_*)\mathbf{C}_*\mathbf{A}(t - t_*)^\top \\ &\quad + \int_{t_*}^{\min(t, t')} \mathbf{A}(t - \tau)\mathbf{L}\sigma^2\mathbf{L}^\top\mathbf{A}(t' - \tau)^\top d\tau \end{aligned}$$

For the practical algorithm we only ever need  $\text{cov}(\mathbf{X}_{t+h}, \mathbf{X}_{t+h})$ , so let

$$\mathbf{Q}(h) = \int_t^{t+h} \mathbf{A}(h)\mathbf{L}\sigma^2\mathbf{L}^\top\mathbf{A}(h)^\top d\tau$$



The advantage of this formulation is that if we view solution as a **filtering problem**, it allows us to express predictive means in **Nordsieck form**.

Assuming:

1. Observations  $\mathbf{z}_t$  are linked to a hidden state  $\mathbf{x}_t$  by a **Linear** operator:

$$\mathbf{z}_t = \mathbf{H} \mathbf{x}_t + \boldsymbol{\xi}_t$$

where  $\boldsymbol{\xi}_t \sim \mathcal{N}(0, R)$

# The Kalman Filter

Assuming:

1. Observations  $\mathbf{z}_t$  are linked to a hidden state  $\mathbf{x}_t$  by a **Linear** operator:

$$\mathbf{z}_t = \mathbf{H} \mathbf{x}_t + \boldsymbol{\xi}_t$$

where  $\boldsymbol{\xi}_t \sim \mathcal{N}(0, R)$

2. We have some **Prediction equation**

$$\mathbf{x}_{t+1}^- = \mathbf{A} \mathbf{x}_t + \boldsymbol{\eta}_t$$

where  $\boldsymbol{\eta}_t \sim \mathcal{N}(0, Q)$

# The Kalman Filter

Assuming:

1. Observations  $\mathbf{z}_t$  are linked to a hidden state  $\mathbf{x}_t$  by a **Linear** operator:

$$\mathbf{z}_t = \mathbf{H} \mathbf{x}_t + \boldsymbol{\xi}_t$$

where  $\boldsymbol{\xi}_t \sim \mathcal{N}(0, R)$

2. We have some **Prediction equation**

$$\mathbf{x}_{t+1}^- = \mathbf{A} \mathbf{x}_t + \boldsymbol{\eta}_t$$

where  $\boldsymbol{\eta}_t \sim \mathcal{N}(0, Q)$

3. We endow  $\mathbf{x}_0$  with a **Gaussian** prior.

# The Kalman Filter

Assuming:

1. Observations  $\mathbf{z}_t$  are linked to a hidden state  $\mathbf{x}_t$  by a **Linear** operator:

$$\mathbf{z}_t = \mathbf{H} \mathbf{x}_t + \boldsymbol{\xi}_t$$

where  $\boldsymbol{\xi}_t \sim \mathcal{N}(0, R)$

2. We have some **Prediction equation**

$$\mathbf{x}_{t+1}^- = \mathbf{A} \mathbf{x}_t + \boldsymbol{\eta}_t$$

where  $\boldsymbol{\eta}_t \sim \mathcal{N}(0, Q)$

3. We endow  $\mathbf{x}_0$  with a **Gaussian** prior.

the Kalman filter describes how to **exactly update** the posterior distribution as new observations are obtained.

More or less...

1. **Predict:**  $\mathbf{x}_{t+1}^- | \mathbf{z}_t$
2. **Update:**  $\mathbf{x}_{t+1} | \mathbf{z}_{t+1}$

Under the assumptions, distributions on  $\mathbf{x}_t$  and  $\mathbf{x}_t^-$  are **Gaussian** for all  $t$ .

More or less...

1. **Predict:**  $\mathbf{x}_{t+1}^- | \mathbf{z}_t$
2. **Update:**  $\mathbf{x}_{t+1} | \mathbf{z}_{t+1}$

Under the assumptions, distributions on  $\mathbf{x}_t$  and  $\mathbf{x}_t^-$  are **Gaussian** for all  $t$ .

Kalman's equations tell us how to efficiently update the mean and covariance matrices for each new piece of information.

Start with a Gaussian distribution  $\mathbf{X}_0$ . For  $n = 0, \dots, N$  do...

1. Compute the predictive distribution  $\mathbf{X}_{t_n}^- | Z_{[n-1]}$
2. Compute  $z_n = f(t_n, \mathbf{X}_{t_n}^-)$  (noiseless)
3. Find  $\mathbf{X}_{t_n} | Z_{[n]}$

Problem: Step 2 is not tractable.



Solution: replace step 2 with something explicitly computable:

1. Compute the predictive distribution  $X_{t_n}^- | Z_{[n-1]}$
2. Compute  $z_n = f(t_n, \mathbb{E}(X_{t_n}^-))$  (noiseless)
3. Find  $X_{t_n} | Z_{[n]}$



## Runge Kutta Means

---

This formulation makes simpler the earlier result<sup>4</sup> regarding Runge Kutta Means.

---

<sup>4</sup>Schober et al. [2014]

This formulation makes simpler the earlier result<sup>4</sup> regarding Runge Kutta Means. In what follows we will use an **Integrated Wiener Process** of order  $q$  (IWP( $q$ )) form for  $X_t$ :

$$dX_t = U_{q+1}Xdt + e_q dW$$

Where  $U_q$  is the **upper shift matrix** of size  $q$ . This gives a convenient closed-form for  $A(h)$ ,  $Q(h)$ .

---

<sup>4</sup>Schober et al. [2014]

This formulation makes simpler the earlier result<sup>4</sup> regarding Runge Kutta Means. In what follows we will use an **Integrated Wiener Process** of order  $q$  (IWP( $q$ )) form for  $\mathbf{X}_t$ :

$$d\mathbf{X}_t = \mathbf{U}_{q+1}\mathbf{X}dt + \mathbf{e}_q dW$$

Where  $\mathbf{U}_q$  is the **upper shift matrix** of size  $q$ . This gives a convenient closed-form for  $\mathbf{A}(h)$ ,  $\mathbf{Q}(h)$ .

Choose the prior mean to be  $\mathbf{m}_{t-1}^- \equiv 0$  and the prior covariance to be  $\mathbf{C}_{t-1}^-$ .

---

<sup>4</sup>Schober et al. [2014]

The predictive mean at  $t = 1$  is equivalent to Explicit Euler.

The **predictive mean** at  $t = 1$  is equivalent to **Explicit Euler**.

$$\begin{array}{ccc} \bullet & \circ & \circ \\ t_{-1} & t_0 & t_1 \end{array}$$

$$\mathbf{m}_{t_{-1}} = \begin{pmatrix} y_0 \\ m_{t_0,1}^- \end{pmatrix}$$
$$\mathbf{C}_{t_{-1}} = \begin{pmatrix} 0 & 0 \\ 0 & C_{t_0,11}^- \end{pmatrix}$$

for some  $m_{t_0,1}^-$  and  $C_{t_0,11}^-$



The predictive mean at  $t = 1$  is equivalent to Explicit Euler.

$t_{-1}$        $t_0$        $t_1$

$$m_{t_0} = \begin{pmatrix} y_0 \\ z_0 \end{pmatrix}$$

$$c_{t_0} = 0$$

The **predictive mean** at  $t = 1$  is equivalent to **Explicit Euler**.



$$m_{t_1}^- = \begin{pmatrix} y_0 + hz_0 \\ z_0 \end{pmatrix}$$

$$C_{t_1}^- = Q(h)$$

The **predictive mean** at  $t = 1$  is equivalent to **Explicit Euler**.

$$\bullet \quad \bullet \quad \bullet$$

$t_{-1} \quad t_0 \quad t_1$

$$m_{t_1} = \begin{pmatrix} y_0 + \frac{h}{2}(z_0 + z_1) \\ z_1 \end{pmatrix}$$

$$C_{t_1} = \sigma^2 \begin{pmatrix} \frac{h^3}{12} & 0 \\ 0 & 0 \end{pmatrix}$$

This corresponds to an RK2 scheme. (Heun's Method,  $\alpha = 1$ ).

Continuing this past  $t_1$  we see that the scheme **no longer matches** any known numerical method (not even Heun's method).

For an IWP(2) model we fix the prior covariance to be  $\mathbf{C}_{t-1}^- = \mathbf{Q}(\tau)$  for some  $\tau > 0$ .

For an IWP(2) model we fix the prior covariance to be  $\mathbf{C}_{t-1}^- = \mathbf{Q}(\tau)$  for some  $\tau > 0$ .



$$\left[ \mathbf{m}_{t_0+h\alpha}^- \right]_0 = y_0 + h\alpha z_0 + \frac{h^2\alpha^2}{2} \left( \frac{4z_0}{\tau} - \frac{20y_0}{3\tau} \right)$$

For an IWP(2) model we fix the prior covariance to be  $\mathbf{C}_{t_{-1}}^- = \mathbf{Q}(\tau)$  for some  $\tau > 0$ .

$$\begin{array}{ccccccc} \bullet & & \bullet & & \bullet & & \circ \\ t_{-1} & & t_0 & t_0 + h\alpha & t_1 & & \end{array}$$

$$\left[ m_{t_0+h\alpha}^- \right]_0 = y_0 + h\alpha z_0 + \frac{h^2\alpha^2}{2} \left( \frac{4z_0}{\tau} - \frac{20y_0}{3\tau} \right)$$

For equivalence with RK2 we require

$$\left[ m_{t_0+h\alpha}^- \right]_0 = y_0 + h\alpha z_0$$

Solution (?): Send  $\tau \rightarrow \infty$

$\tau \rightarrow \infty$

- Equivalent to using an **improper prior**.
- More severe ramifications for continuation.
- Challenging computationally.



$\tau \rightarrow \infty$

- Equivalent to using an **improper prior**.
- More severe ramifications for continuation.
- Challenging computationally.

Difficult to justify given that the **domain** is restricted to  $[t_0, T]$

# Convergence Analysis and Nordiseck Methods

---

The posterior mean **does not** correspond to an RK method past  $t_1$ ...

⇒ no global error results from this route!

The posterior mean **does not** correspond to an RK method past  $t_1$ ...

⇒ no global error results from this route!

It **does** correspond to a **general linear method** in Nordsieck form.

We re-scale the state vector:

$$\tilde{\mathbf{X}} = \begin{pmatrix} Y_t \\ hY'_t \\ \frac{h^2}{2!}Y''_t \\ \vdots \\ \frac{h^q}{q!}Y_t^{(q)} \end{pmatrix} = \mathbf{B}\mathbf{X}_t$$

yielding the new SDE

$$d\tilde{\mathbf{X}}_t = \mathbf{B}\mathbf{U}_{q+1}\mathbf{B}^{-1}\tilde{\mathbf{X}}_tdt + \mathbf{B}\mathbf{e}_qdW$$

which makes  $\tilde{\mathbf{A}}(h)$  independent of  $h$ .

**Proposition 1**

*The probabilistic Nordsieck method arising from the once-integrated Wiener process is equivalent in predictive posterior mean with the trapezoidal rule.*

### Proposition 1

*The probabilistic Nordsieck method arising from the once-integrated Wiener process is equivalent in predictive posterior mean with the trapezoidal rule.*

This means that the IWP(1) probabilistic Nordsieck method is **actually of order 2 rather than order 1!**

## Theorem 1

*The predictive posterior mean of the IWP(2) with fixed step size  $h$  is a third order Nordsieck method.*



# Calibration and Experiments

---

Three parameters to tune:  $q$ ,  $\sigma^2$ ,  $h$ .  $q$  is fixed to 2.

Three parameters to tune:  $q$ ,  $\sigma^2$ ,  $h$ .  $q$  is fixed to 2.

As with standard solvers we **adapt**  $h_n$  at each step to not exceed a tolerance.

Three parameters to tune:  $q, \sigma^2, h$ .  $q$  is fixed to 2.

As with standard solvers we **adapt**  $h_n$  at each step to not exceed a tolerance.

$\sigma$  is chosen by maximising the marginal likelihood of the **residual**:

$$\hat{\sigma} := \arg \min_{\sigma} p \left( z_n - [m_{t_n}^-]_1 \mid \sigma \right)$$

[To the paper!]

## Conclusions

---

- The choice of  $z_n = f(t_n, \mathbb{E}(X_{t_n}^-))$ .

- The choice of  $z_n = f(t_n, \mathbb{E}(X_{t_n}^-))$ .
- The filtration assumption.



- How does this relate to the fully Bayesian posterior distribution?

- How does this relate to the fully Bayesian posterior distribution?
- What about **implicit schemes**?

Thanks!

# References

---

Michael Schober, David K Duvenaud, and Philipp Hennig. Probabilistic ode solvers with runge-kutta means. In *Advances in neural information processing systems*, pages 739–747, 2014.

Michael Schober, Simo Särkkä, and Philipp Hennig. A probabilistic model for the numerical solution of initial value problems. *arXiv preprint arXiv:1610.05261*, 2016.

Robert D Skeel. Equivalent forms of multistep formulas. *Mathematics of Computation*, 33(148):1229–1250, 1979.

Onur Teymur, Kostas Zygalakis, and Ben Calderhead. Probabilistic linear multistep methods. In *Advances in Neural Information Processing Systems*, pages 4314–4321, 2016.