

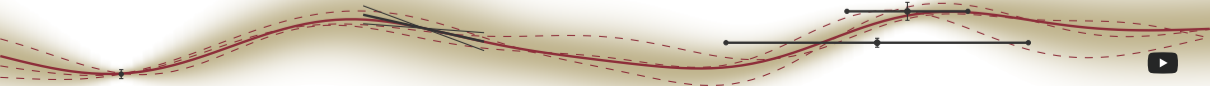
NUMERICS OF MACHINE LEARNING
LECTURE 07
PROBABILISTIC NUMERICAL ODE SOLVERS

Nathanael Bosch & Jonathan Schmidt
1 December 2022

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



FACULTY OF SCIENCE
DEPARTMENT OF COMPUTER SCIENCE
CHAIR FOR THE METHODS OF MACHINE LEARNING



Two weeks ago: **State-space models and extended Kalman filters/smoothers**

- ▶ “How to estimate the *state* of a dynamical system from *observations*”

Last week: **Ordinary differential equations and how to solve them**

- ▶ “How to *simulate*, approximately, a deterministic dynamical system”

This week: **ODE simulation as probabilistic inference**

- ▶ “How to treat ODEs as the state estimation problem that they really are”



Two weeks ago: **State-space models and extended Kalman filters/smoothers**

- ▶ “How to estimate the *state* of a dynamical system from *observations*”

Last week: **Ordinary differential equations and how to solve them**

- ▶ “How to *simulate*, approximately, a deterministic dynamical system”

This week: **ODE simulation as probabilistic inference**

- ▶ “How to treat ODEs as the state estimation problem that they really are”
 \Rightarrow ***Probabilistic* numerical ODE solvers**



Recap: Numerical Ordinary Differential Equation Solvers



Recap: Ordinary Differential Equations and how to solve them

Numerical ODE solvers try to estimate an unknown function by evaluating the vector field

$$\dot{x}(t) = f(x(t), t)$$

with $t \in [0, T]$, vector field $f: \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^d$, and initial value $x(0) = x_0$. Goal: "Find x ".

Recap: Ordinary Differential Equations and how to solve them

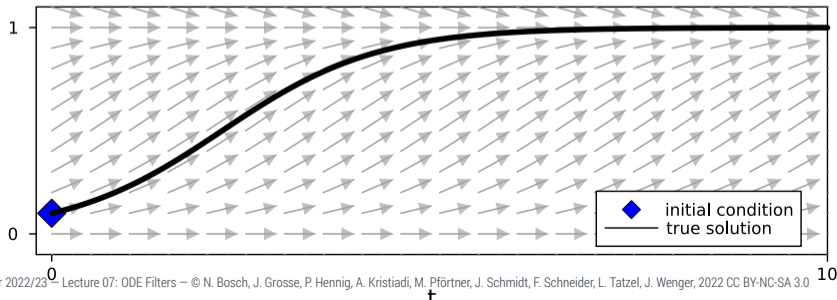
Numerical ODE solvers try to estimate an unknown function by evaluating the vector field

$$\dot{x}(t) = f(x(t), t)$$

with $t \in [0, T]$, vector field $f: \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^d$, and initial value $x(0) = x_0$. Goal: "Find x ".

► Simple example: Logistic ODE

$$\dot{x}(t) = x(t)(1 - x(t)), \quad t \in [0, 10], \quad x(0) = 0.1.$$



Numerical ODE solvers try to estimate an unknown function by evaluating the vector field

$$\dot{x}(t) = f(x(t), t)$$

with $t \in [0, T]$, vector field $f: \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^d$, and initial value $x(0) = x_0$. Goal: "Find x ".

Numerical ODE solvers:

- ▶ Forward Euler:

$$\hat{x}(t+h) = \hat{x}(t) + h \cdot f(\hat{x}(t), t)$$

- ▶ Backward Euler:

$$\hat{x}(t+h) = \hat{x}(t) + h \cdot f(\hat{x}(t+h), t+h)$$

- ▶ Runge-Kutta:

$$\hat{x}(t+h) = \hat{x}(t) + h \cdot \sum_{i=1}^s b_i f(\tilde{x}_i, t + c_i h)$$

Recap: Ordinary Differential Equations and how to solve them

Numerical ODE solvers try to estimate an unknown function by evaluating the vector field

$$\dot{x}(t) = f(x(t), t)$$

with $t \in [0, T]$, vector field $f: \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^d$, and initial value $x(0) = x_0$. Goal: "Find x ".

Numerical ODE solvers:

- ▶ Forward Euler:

$$\hat{x}(t+h) = \hat{x}(t) + h \cdot f(\hat{x}(t), t)$$

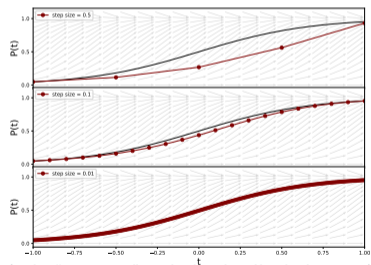
- ▶ Backward Euler:

$$\hat{x}(t+h) = \hat{x}(t) + h \cdot f(\hat{x}(t+h), t+h)$$

- ▶ Runge–Kutta:

$$\hat{x}(t+h) = \hat{x}(t) + h \cdot \sum_{i=1}^S b_i f(\tilde{x}_i, t + c_i h)$$

Forward Euler for different step sizes:



(It is "correct" only in the limit $h \rightarrow 0$!)

Recap: Ordinary Differential Equations and how to solve them

Numerical ODE solvers try to estimate an unknown function by evaluating the vector field

$$\dot{x}(t) = f(x(t), t)$$

with $t \in [0, T]$, vector field $f: \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}^d$, and initial value $x(0) = x_0$. Goal: "Find x ".

Numerical ODE solvers:

- ▶ Forward Euler:

$$\hat{x}(t+h) = \hat{x}(t) + h \cdot f(\hat{x}(t), t)$$

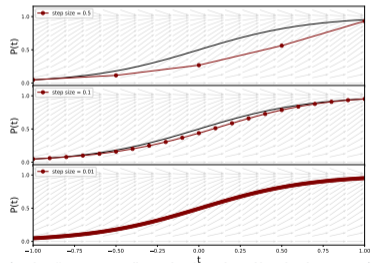
- ▶ Backward Euler:

$$\hat{x}(t+h) = \hat{x}(t) + h \cdot f(\hat{x}(t+h), t+h)$$

- ▶ Runge-Kutta:

$$\hat{x}(t+h) = \hat{x}(t) + h \cdot \sum_{i=1}^S b_i f(\tilde{x}_i, t + c_i h)$$

Forward Euler for different step sizes:



(It is "correct" only in the limit $h \rightarrow 0$!)

Numerical ODE solvers **estimate** $x(t)$ by evaluating f on a discrete set of points.

Recap: Bayesian State Estimation with Extended Kalman filtering and smoothing



Non-linear Gaussian state-estimation problem:

Initial distribution: $x_0 \sim \mathcal{N}(\mu_0, \Sigma_0),$

Prior / dynamics model: $x_{i+1} | x_i \sim \mathcal{N}(f(x_i), Q_i),$

Likelihood / measurement model: $y_i | x_i \sim \mathcal{N}(h(x_i), R_i),$

Data: $\mathcal{D} = \{y_i\}_{i=1}^N.$

Non-linear Gaussian state-estimation problem:

Initial distribution: $x_0 \sim \mathcal{N}(\mu_0, \Sigma_0),$

Prior / dynamics model: $x_{i+1} | x_i \sim \mathcal{N}(f(x_i), Q_i),$

Likelihood / measurement model: $y_i | x_i \sim \mathcal{N}(h(x_i), R_i),$

Data: $\mathcal{D} = \{y_i\}_{i=1}^N.$

EKF/EKS: The extended Kalman filter/smoothener recursively computes Gaussian approximations:

Predict: $p(x_i | y_{1:i-1}) \approx \mathcal{N}(x_i; \mu_i^P, \Sigma_i^P),$

Filter: $p(x_i | y_{1:i}) \approx \mathcal{N}(x_i; \mu_i, \Sigma_i),$

Smooth: $p(x_i | y_{1:N}) \approx \mathcal{N}(x_i; \mu_i^S, \Sigma_i^S),$

Likelihood: $p(y_i | y_{1:i-1}) \approx \mathcal{N}(y_i; \hat{y}_i, S_i).$

Recap: Extended Kalman filtering and smoothing

EKF/EKS as introduced in lecture 5

Non-linear Gaussian state-estimation problem:

Initial distribution: $x_0 \sim \mathcal{N}(\mu_0, \Sigma_0)$,

Prior / dynamics model: $x_{i+1} | x_i \sim \mathcal{N}(f(x_i), Q_i)$,

Likelihood / measurement model: $y_i | x_i \sim \mathcal{N}(h(x_i), R_i)$,

Data: $\mathcal{D} = \{y_i\}_{i=1}^N$.

PREDICT

$$\mu_{i+1}^P = f(\mu_i),$$

$$\Sigma_{i+1}^P = J_f(\mu_i)\Sigma_i J_f(\mu_i)^\top + Q_i.$$

EKF/EKS: The extended Kalman filter/smoothener recursively computes Gaussian approximations:

Predict: $p(x_i | y_{1:i-1}) \approx \mathcal{N}(x_i; \mu_i^P, \Sigma_i^P)$,

Filter: $p(x_i | y_{1:i}) \approx \mathcal{N}(x_i; \mu_i, \Sigma_i)$,

Smooth: $p(x_i | y_{1:N}) \approx \mathcal{N}(x_i; \mu_i^S, \Sigma_i^S)$,

Likelihood: $p(y_i | y_{1:i-1}) \approx \mathcal{N}(y_i; \hat{y}_i, S_i)$.

Recap: Extended Kalman filtering and smoothing

EKF/EKS as introduced in lecture 5

Non-linear Gaussian state-estimation problem:

Initial distribution: $x_0 \sim \mathcal{N}(\mu_0, \Sigma_0),$

Prior / dynamics model: $x_{i+1} | x_i \sim \mathcal{N}(f(x_i), Q_i),$

Likelihood / measurement model: $y_i | x_i \sim \mathcal{N}(h(x_i), R_i),$

Data: $\mathcal{D} = \{y_i\}_{i=1}^N.$

EKF/EKS: The extended Kalman filter/smoothener recursively computes Gaussian approximations:

Predict: $p(x_i | y_{1:i-1}) \approx \mathcal{N}(x_i; \mu_i^P, \Sigma_i^P),$

Filter: $p(x_i | y_{1:i}) \approx \mathcal{N}(x_i; \mu_i, \Sigma_i),$

Smooth: $p(x_i | y_{1:N}) \approx \mathcal{N}(x_i; \mu_i^S, \Sigma_i^S),$

Likelihood: $p(y_i | y_{1:i-1}) \approx \mathcal{N}(y_i; \hat{y}_i, S_i).$

PREDICT

$$\mu_{i+1}^P = f(\mu_i),$$

$$\Sigma_{i+1}^P = J_f(\mu_i) \Sigma_i J_f(\mu_i)^\top + Q_i.$$

UPDATE

$$\hat{z}_i = h(\mu_i^P),$$

$$S_i = J_h(\mu_i^P) \Sigma_i^P J_h(\mu_i^P)^\top + R_i,$$

$$K_i = \Sigma_i^P J_h(\mu_i^P)^\top S_i^{-1},$$

$$\mu_i = \mu_i^P + K_i (z_i - \hat{z}_i),$$

$$\Sigma_i = \Sigma_i^P - K_i S_i K_i^\top.$$

Recap: Extended Kalman filtering and smoothing

EKF/EKS as introduced in lecture 5

Non-linear Gaussian state-estimation problem:

Initial distribution: $x_0 \sim \mathcal{N}(\mu_0, \Sigma_0)$,

Prior / dynamics model: $x_{i+1} | x_i \sim \mathcal{N}(f(x_i), Q_i)$,

Likelihood / measurement model: $y_i | x_i \sim \mathcal{N}(h(x_i), R_i)$,

Data: $\mathcal{D} = \{y_i\}_{i=1}^N$.

EKF/EKS: The extended Kalman filter/smoothener recursively computes Gaussian approximations:

Predict: $p(x_i | y_{1:i-1}) \approx \mathcal{N}(x_i; \mu_i^P, \Sigma_i^P)$,

Filter: $p(x_i | y_{1:i}) \approx \mathcal{N}(x_i; \mu_i, \Sigma_i)$,

Smooth: $p(x_i | y_{1:N}) \approx \mathcal{N}(x_i; \mu_i^S, \Sigma_i^S)$,

Likelihood: $p(y_i | y_{1:i-1}) \approx \mathcal{N}(y_i; \hat{y}_i, S_i)$.

PREDICT

$$\mu_{i+1}^P = f(\mu_i),$$

$$\Sigma_{i+1}^P = J_f(\mu_i) \Sigma_i J_f(\mu_i)^\top + Q_i.$$

UPDATE

$$\hat{z}_i = h(\mu_i^P),$$

$$S_i = J_h(\mu_i^P) \Sigma_i^P J_h(\mu_i^P)^\top + R_i,$$

$$K_i = \Sigma_i^P J_h(\mu_i^P)^\top S_i^{-1},$$

$$\mu_i = \mu_i^P + K_i (z_i - \hat{z}_i),$$

$$\Sigma_i = \Sigma_i^P - K_i S_i K_i^\top.$$

SMOOTH: See lecture 5.

Today: *Probabilistic* numerical ODE solutions

or “how to treat ODEs as the state estimation problem that they really are”



Probabilistic numerical ODE solutions

How to treat ODEs as the state estimation problem that they really are



Probabilistic numerical ODE solutions

How to treat ODEs as the state estimation problem that they really are

$$p \left(x(t) \mid x(0) = x_0, \{\dot{x}(t_n) = f(x(t_n), t_n)\}_{n=1}^N \right)$$



Probabilistic numerical ODE solutions

How to treat ODEs as the state estimation problem that they really are

$$p \left(x(t) \mid x(0) = x_0, \{ \dot{x}(t_n) = f(x(t_n), t_n) \}_{n=1}^N \right)$$

We want *fast* (approximate) inference



Probabilistic numerical ODE solutions

How to treat ODEs as the state estimation problem that they really are

$$p \left(x(t) \mid x(0) = x_0, \{\dot{x}(t_n) = f(x(t_n), t_n)\}_{n=1}^N \right)$$

We want *fast* (approximate) inference

⇒ Gaussian filtering and smoothing (it's $\mathcal{O}(N)$!)



Probabilistic numerical ODE solutions

How to treat ODEs as the state estimation problem that they really are

$$p \left(x(t) \mid x(0) = x_0, \{ \dot{x}(t_n) = f(x(t_n), t_n) \}_{n=1}^N \right)$$

We want *fast* (approximate) inference

⇒ Gaussian filtering and smoothing (it's $\mathcal{O}(N)$!)

⇒ We need to construct a state-space model:

1. Prior:
2. Likelihood:
3. Data:



Probabilistic numerical ODE solutions

How to treat ODEs as the state estimation problem that they really are

$$p \left(x(t) \mid x(0) = x_0, \{ \dot{x}(t_n) = f(x(t_n), t_n) \}_{n=1}^N \right)$$

We want *fast* (approximate) inference

⇒ Gaussian filtering and smoothing (it's $\mathcal{O}(N)$!)

⇒ We need to construct a state-space model:

1. **Prior:**
2. Likelihood:
3. Data:



The Prior: Describing the evolution of the “state”

What is the “state”, and how can we model it?



The Prior: Describing the evolution of the “state”

What is the “state”, and how can we model it?

- ▶ **What is the “state”?** Only $x(t)$ is not sufficient to fully describe the dynamical system.

The Prior: Describing the evolution of the “state”

What is the “state”, and how can we model it?

- ▶ **What is the “state”?** Only $x(t)$ is not sufficient to fully describe the dynamical system.
- ▶ **Motivation** from Taylor series expansions:

$$x(t+h) = x(t) + h\dot{x}(t) + \frac{h^2}{2}\ddot{x}(t) + \dots = \sum_{k=0}^{\infty} \frac{h^k}{k!} x^{(k)}(t).$$

⇒ Having access to *all* derivatives *would* fully describe the dynamical system.

The Prior: Describing the evolution of the “state”

What is the “state”, and how can we model it?

- ▶ **What is the “state”?** Only $x(t)$ is not sufficient to fully describe the dynamical system.
- ▶ **Motivation** from Taylor series expansions:

$$x(t+h) = x(t) + h\dot{x}(t) + \frac{h^2}{2}\ddot{x}(t) + \dots = \sum_{k=0}^{\infty} \frac{h^k}{k!} x^{(k)}(t).$$

⇒ Having access to *all* derivatives *would* fully describe the dynamical system.

- ▶ **Truncating the Taylor series expansion** greatly simplifies things, but introduces an error term:

$$x(t+h) = \sum_{k=0}^q \frac{h^k}{k!} x^{(k)}(t) + \mathcal{O}(h^{q+1}).$$

The Prior: Describing the evolution of the “state”

What is the “state”, and how can we model it?

- ▶ **What is the “state”?** Only $x(t)$ is not sufficient to fully describe the dynamical system.
- ▶ **Motivation** from Taylor series expansions:

$$x(t+h) = x(t) + h\dot{x}(t) + \frac{h^2}{2}\ddot{x}(t) + \dots = \sum_{k=0}^{\infty} \frac{h^k}{k!} x^{(k)}(t).$$

⇒ Having access to *all* derivatives *would* fully describe the dynamical system.

- ▶ **Truncating the Taylor series expansion** greatly simplifies things, but introduces an error term:

$$x(t+h) = \sum_{k=0}^q \frac{h^k}{k!} x^{(k)}(t) + \mathcal{O}(h^{q+1}).$$

- ▶ **Towards a “prior process”:** Modeling the error with a random variable gives

$$x(t+h) = \sum_{k=0}^q \frac{h^k}{k!} x^{(k)}(t) + \epsilon(t,h), \quad \epsilon(t,h) \sim \mathcal{N}\left(0, \frac{h^{2q+1}}{(2q+1)(q!)^2}\right)$$

The Prior: Describing the evolution of the “state”

What is the “state”, and how can we model it?

- ▶ **What is the “state”?** Only $x(t)$ is not sufficient to fully describe the dynamical system.
- ▶ **Motivation** from Taylor series expansions:

$$x(t+h) = x(t) + h\dot{x}(t) + \frac{h^2}{2}\ddot{x}(t) + \dots = \sum_{k=0}^{\infty} \frac{h^k}{k!} x^{(k)}(t).$$

⇒ Having access to *all* derivatives *would* fully describe the dynamical system.

- ▶ **Truncating the Taylor series expansion** greatly simplifies things, but introduces an error term:

$$x(t+h) = \sum_{k=0}^q \frac{h^k}{k!} x^{(k)}(t) + \mathcal{O}(h^{q+1}).$$

- ▶ **Towards a “prior process”:** Modeling the error with a random variable gives

$$x(t+h) = \sum_{k=0}^q \frac{h^k}{k!} x^{(k)}(t) + \epsilon(t,h), \quad \epsilon(t,h) \sim \mathcal{N}\left(0, \frac{h^{2q+1}}{(2q+1)(q!)^2}\right)$$

- ▶ **This motivates a “state”:** $X(t) = [x(t), \dot{x}(t), \ddot{x}(t), \dots, x^{(q)}(t)]^\top$. *Details on the next slide.*

The Prior: Describing the evolution of the “state”

The q -times integrated Wiener process

The prior model: q -times integrated Wiener process (IWP(q)).

The Prior: Describing the evolution of the “state”

The q -times integrated Wiener process

The prior model: q -times integrated Wiener process (IWP(q)).

Let $X(t)$ be of the form

$$X(t) = \left[X^{(0)}(t), X^{(1)}(t), \dots, X^{(q)}(t) \right]^T,$$

chosen such that $X^{(i)}(t)$ models the i -th derivative of $x(t)$.

The Prior: Describing the evolution of the “state”

The q -times integrated Wiener process

Formulas in: *Kersting et al, “Convergence Rates of Gaussian ODE Filters”, 2020*

The prior model: q -times integrated Wiener process (IWP(q)).

Let $X(t)$ be of the form

$$X(t) = \left[X^{(0)}(t), X^{(1)}(t), \dots, X^{(q)}(t) \right]^T,$$

chosen such that $X^{(i)}(t)$ models the i -th derivative of $x(t)$.

The IWP(q) has known discrete-time transitions:

$$X(t+h) \mid X(t) \sim \mathcal{N} \left(X(t+h); A(h)X(t), \sigma^2 Q(h) \right),$$

$$[A(h)]_{ij} = \mathbb{1}_{i \leq j} \frac{h^{j-i}}{(j-i)!},$$

$$[Q(h)]_{ij} = \frac{h^{2q+1-i-j}}{(2q+1-i-j)(q-i)!(q-j)!},$$

for any $i, j = 0, \dots, q$.

The Prior: Describing the evolution of the “state”

The q -times integrated Wiener process

Formulas in: *Kersting et al, “Convergence Rates of Gaussian ODE Filters”, 2020*

The prior model: q -times integrated Wiener process (IWP(q)).

Let $X(t)$ be of the form

$$X(t) = \left[X^{(0)}(t), X^{(1)}(t), \dots, X^{(q)}(t) \right]^T,$$

chosen such that $X^{(i)}(t)$ models the i -th derivative of $x(t)$.

The IWP(q) has known discrete-time transitions:

$$X(t+h) \mid X(t) \sim \mathcal{N} \left(X(t+h); A(h)X(t), \sigma^2 Q(h) \right),$$

$$[A(h)]_{ij} = \mathbb{1}_{i \leq j} \frac{h^{j-i}}{(j-i)!},$$

$$[Q(h)]_{ij} = \frac{h^{2q+1-i-j}}{(2q+1-i-j)(q-i)!(q-j)!},$$

for any $i, j = 0, \dots, q$.

Example: IWP(2)

$$A(h) = \begin{pmatrix} 1 & h & \frac{h^2}{2} \\ 0 & 1 & h \\ 0 & 0 & 1 \end{pmatrix},$$

$$Q(h) = \begin{pmatrix} \frac{h^5}{20} & \frac{h^4}{8} & \frac{h^3}{6} \\ \frac{h^4}{8} & \frac{h^3}{3} & \frac{h^2}{2} \\ \frac{h^3}{6} & \frac{h^2}{2} & h \end{pmatrix}.$$



The Prior: Describing the evolution of the “state”

The q -times integrated Wiener process

Formulas in: Kersting et al, “Convergence Rates of Gaussian ODE Filters”, 2020

The prior model: q -times integrated Wiener process (IWP(q)).

Let $X(t)$ be of the form

$$X(t) = \left[X^{(0)}(t), X^{(1)}(t), \dots, X^{(q)}(t) \right]^T,$$

chosen such that $X^{(i)}(t)$ models the i -th derivative of $x(t)$.

The IWP(q) has known discrete-time transitions:

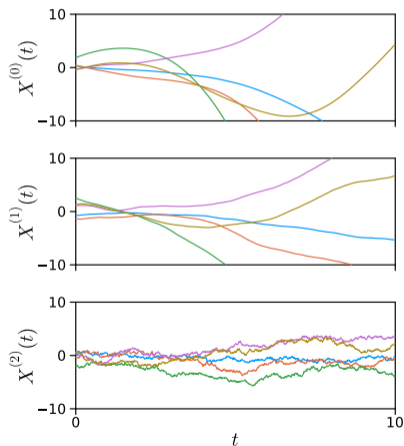
$$X(t+h) \mid X(t) \sim \mathcal{N} \left(X(t+h); A(h)X(t), \sigma^2 Q(h) \right),$$

$$[A(h)]_{ij} = \mathbb{1}_{i \leq j} \frac{h^{j-i}}{(j-i)!},$$

$$[Q(h)]_{ij} = \frac{h^{2q+1-i-j}}{(2q+1-i-j)(q-i)!(q-j)!},$$

for any $i, j = 0, \dots, q$.

Example: IWP(2)





The Prior: Describing the evolution of the “state”

The q -times integrated Wiener process

Formulas in: Kersting et al, “Convergence Rates of Gaussian ODE Filters”, 2020

The prior model: q -times integrated Wiener process (IWP(q)).

Let $X(t)$ be of the form

$$X(t) = \left[X^{(0)}(t), X^{(1)}(t), \dots, X^{(q)}(t) \right]^T,$$

chosen such that $X^{(i)}(t)$ models the i -th derivative of $x(t)$.

The IWP(q) has known discrete-time transitions:

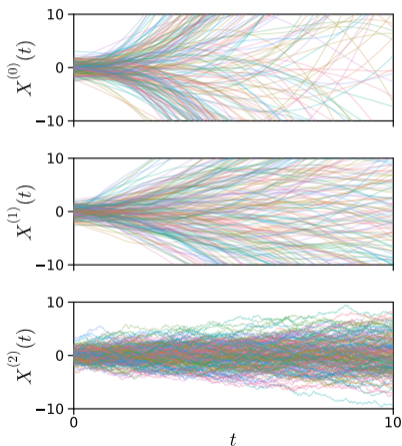
$$X(t+h) \mid X(t) \sim \mathcal{N} \left(X(t+h); A(h)X(t), \sigma^2 Q(h) \right),$$

$$[A(h)]_{ij} = \mathbb{1}_{i \leq j} \frac{h^{j-i}}{(j-i)!},$$

$$[Q(h)]_{ij} = \frac{h^{2q+1-i-j}}{(2q+1-i-j)(q-i)!(q-j)!},$$

for any $i, j = 0, \dots, q$.

Example: IWP(2)



The Prior: Describing the evolution of the “state”

The q -times integrated Wiener process

Formulas in: Kersting et al, “Convergence Rates of Gaussian ODE Filters”, 2020

The prior model: q -times integrated Wiener process (IWP(q)).

Let $X(t)$ be of the form

$$X(t) = \left[X^{(0)}(t), X^{(1)}(t), \dots, X^{(q)}(t) \right]^T,$$

chosen such that $X^{(i)}(t)$ models the i -th derivative of $x(t)$.

The IWP(q) has known discrete-time transitions:

$$X(t+h) | X(t) \sim \mathcal{N} \left(X(t+h); A(h)X(t), \sigma^2 Q(h) \right),$$

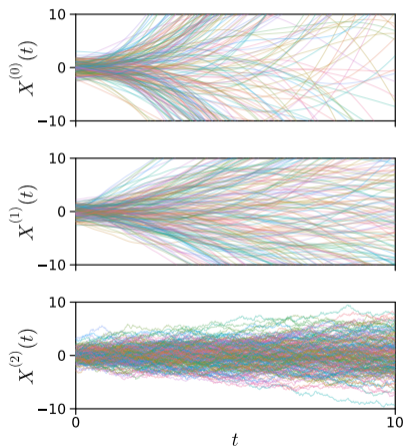
$$[A(h)]_{ij} = \mathbb{1}_{i \leq j} \frac{h^{j-i}}{(j-i)!},$$

$$[Q(h)]_{ij} = \frac{h^{2q+1-i-j}}{(2q+1-i-j)(q-i)!(q-j)!},$$

for any $i, j = 0, \dots, q$.

For later convenience: Define projection matrices $E_i X = X^{(i)}$.

Example: IWP(2)



How to treat ODEs as the state estimation problem that they really are

$$p \left(x(t) \mid x(0) = x_0, \{ \dot{x}(t_n) = f(x(t_n), t_n) \}_{n=1}^N \right)$$

We want *fast* (approximate) inference \Rightarrow Gaussian filtering and smoothing (it's $\mathcal{O}(N)$!)

1. **Prior:** q -times integrated Wiener process prior:

$$X(t+h) \mid X(t) \sim \mathcal{N}(X(t+h); A(h)X(t), Q(h))$$

2. Likelihood:
3. Data:



Probabilistic numerical ODE solutions

How to treat ODEs as the state estimation problem that they really are

$$p \left(x(t) \mid x(0) = x_0, \{ \dot{x}(t_n) = f(x(t_n), t_n) \}_{n=1}^N \right)$$

We want *fast* (approximate) inference \Rightarrow Gaussian filtering and smoothing (it's $\mathcal{O}(N)$!)

1. Prior: q -times integrated Wiener process prior:

$$X(t+h) \mid X(t) \sim \mathcal{N}(X(t+h); A(h)X(t), Q(h))$$

2. Likelihood:
3. Data:



The likelihood model and the data – aka. “*The information operator*”

The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- ▶ **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

The likelihood model and the data – aka. “*The information operator*”

The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

The likelihood model and the data – aka. “*The information operator*”

The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

$$0 = X^{(1)}(t) - f\left(X^{(0)}(t), t\right)$$

The likelihood model and the data – aka. “*The information operator*”

The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

$$0 = X^{(1)}(t) - f\left(X^{(0)}(t), t\right) =: m(X(t), t).$$

The likelihood model and the data – aka. “*The information operator*”

The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

$$0 = X^{(1)}(t) - f\left(X^{(0)}(t), t\right) =: m(X(t), t).$$

- **Easier goal:** Satisfy the ODE *on a discrete time grid* $\{t_i\}_{i=1}^N$

$$\dot{x}(t_i) = f(x(t_i), t_i), \quad i = 1, \dots, N.$$

The likelihood model and the data – aka. “*The information operator*”

The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

$$0 = X^{(1)}(t) - f\left(X^{(0)}(t), t\right) =: m(X(t), t).$$

- **Easier goal:** Satisfy the ODE *on a discrete time grid* $\{t_i\}_{i=1}^N$

$$\dot{x}(t_i) = f(x(t_i), t_i), \quad i = 1, \dots, N.$$

$$\Leftrightarrow m(X(t_i), t_i) = 0$$

The likelihood model and the data – aka. “*The information operator*”

The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- ▶ **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

$$0 = X^{(1)}(t) - f\left(X^{(0)}(t), t\right) =: m(X(t), t).$$

- ▶ **Easier goal:** Satisfy the ODE *on a discrete time grid* $\{t_i\}_{i=1}^N$

$$\dot{x}(t_i) = f(x(t_i), t_i), \quad i = 1, \dots, N.$$

$$\Leftrightarrow m(X(t_i), t_i) = 0$$

- ▶ This motivates the **likelihood model** and **data**:

$$Z(t_i) \mid X(t_i) \sim \mathcal{N}(m(X(t_i), t_i), R)$$

$$z_i \triangleq 0, \quad i = 1, \dots, N.$$

where z_i is a realization of $Z(t_i)$.

The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- ▶ **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

$$0 = X^{(1)}(t) - f\left(X^{(0)}(t), t\right) =: m(X(t), t).$$

- ▶ **Easier goal:** Satisfy the ODE *on a discrete time grid* $\{t_i\}_{i=1}^N$

$$\dot{x}(t_i) = f(x(t_i), t_i), \quad i = 1, \dots, N.$$

$$\Leftrightarrow m(X(t_i), t_i) = 0$$

- ▶ This motivates the *noiseless likelihood model* and **data**:

$$Z(t_i) \mid X(t_i) \sim \mathcal{N}(m(X(t_i), t_i), 0)$$

$$z_i \triangleq 0, \quad i = 1, \dots, N.$$

where z_j is a realization of $Z(t_j)$.

The likelihood model and the data – aka. “*The information operator*”

The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- ▶ **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

$$0 = X^{(1)}(t) - f\left(X^{(0)}(t), t\right) =: m(X(t), t).$$

- ▶ **Easier goal:** Satisfy the ODE *on a discrete time grid* $\{t_i\}_{i=1}^N$

$$\dot{x}(t_i) = f(x(t_i), t_i), \quad i = 1, \dots, N.$$

$$\Leftrightarrow m(X(t_i), t_i) = 0$$

- ▶ This motivates the *noiseless likelihood model* and **data**:

$$Z(t_i) \mid X(t_i) \sim \delta(m(X(t_i), t_i))$$

$$z_i \triangleq 0, \quad i = 1, \dots, N.$$

where z_i is a realization of $Z(t_i)$.

(δ is the Dirac distribution)

The likelihood model and the data – aka. “The information operator”



The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

$$0 = X^{(1)}(t) - f\left(X^{(0)}(t), t\right) =: m(X(t), t).$$

- **Easier goal:** Satisfy the ODE *on a discrete time grid* $\{t_i\}_{i=1}^N$

$$\dot{x}(t_i) = f(x(t_i), t_i), \quad i = 1, \dots, N.$$

$$\Leftrightarrow m(X(t_i), t_i) = 0$$

- This motivates the *noiseless likelihood model* and **data**:

$$Z(t_i) | X(t_i) \sim \delta(m(X(t_i), t_i))$$

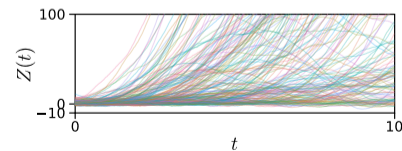
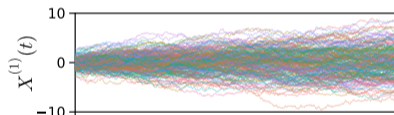
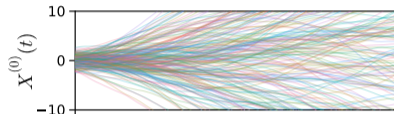
$$z_i \triangleq 0, \quad i = 1, \dots, N.$$

where z_i is a realization of $Z(t_i)$.

(δ is the Dirac distribution)

Example: Logistic ODE $\dot{x} = x(1 - x)$

Prior samples



(here: $Z = X^{(1)} - X^{(0)}(1 - X^{(0)})$)

The likelihood model and the data – aka. “The information operator”



The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

$$0 = X^{(1)}(t) - f\left(X^{(0)}(t), t\right) =: m(X(t), t).$$

- **Easier goal:** Satisfy the ODE *on a discrete time grid* $\{t_i\}_{i=1}^N$

$$\dot{x}(t_i) = f(x(t_i), t_i), \quad i = 1, \dots, N.$$

$$\Leftrightarrow m(X(t_i), t_i) = 0$$

- This motivates the *noiseless likelihood model* and **data**:

$$Z(t_i) | X(t_i) \sim \delta(m(X(t_i), t_i))$$

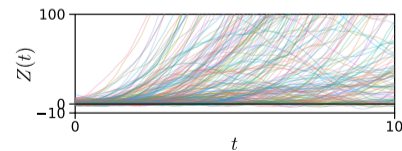
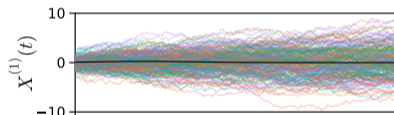
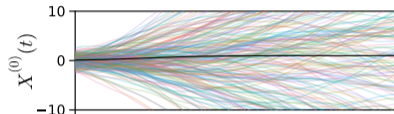
$$z_i \triangleq 0, \quad i = 1, \dots, N.$$

where z_i is a realization of $Z(t_i)$.

(δ is the Dirac distribution)

Example: Logistic ODE $\dot{x} = x(1 - x)$

Prior samples & ODE solution



(here: $Z = X^{(1)} - X^{(0)}(1 - X^{(0)})$)

The likelihood model and the data – aka. “*The information operator*”



The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

$$0 = X^{(1)}(t) - f\left(X^{(0)}(t), t\right) =: m(X(t), t).$$

- **Easier goal:** Satisfy the ODE *on a discrete time grid* $\{t_i\}_{i=1}^N$

$$\dot{x}(t_i) = f(x(t_i), t_i), \quad i = 1, \dots, N.$$

$$\Leftrightarrow m(X(t_i), t_i) = 0$$

- This motivates the *noiseless likelihood model* and **data**:

$$Z(t_i) \mid X(t_i) \sim \delta(m(X(t_i), t_i))$$

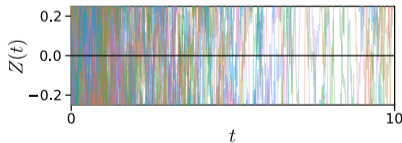
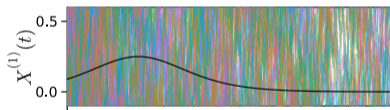
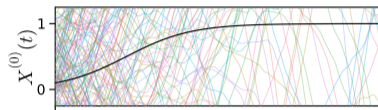
$$z_i \triangleq 0, \quad i = 1, \dots, N.$$

where z_i is a realization of $Z(t_i)$.

(δ is the Dirac distribution)

Example: Logistic ODE $\dot{x} = x(1 - x)$

Prior samples & ODE solution (zoomed)



(here: $Z = X^{(1)} - X^{(0)}(1 - X^{(0)})$)

The likelihood model and the data – aka. “The information operator”



The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

$$0 = X^{(1)}(t) - f\left(X^{(0)}(t), t\right) =: m(X(t), t).$$

- **Easier goal:** Satisfy the ODE *on a discrete time grid* $\{t_i\}_{i=1}^N$

$$\dot{x}(t_i) = f(x(t_i), t_i), \quad i = 1, \dots, N.$$

$$\Leftrightarrow m(X(t_i), t_i) = 0$$

- This motivates the *noiseless likelihood model* and **data**:

$$Z(t_i) | X(t_i) \sim \delta(m(X(t_i), t_i))$$

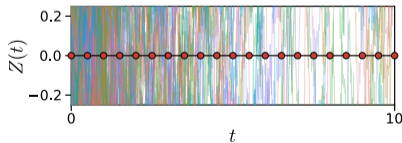
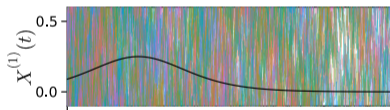
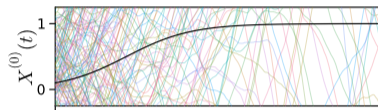
$$z_i \triangleq 0, \quad i = 1, \dots, N.$$

where z_i is a realization of $Z(t_i)$.

(δ is the Dirac distribution)

Example: Logistic ODE $\dot{x} = x(1 - x)$

Prior samples & ODE solution & “Data”



(here: $Z = X^{(1)} - X^{(0)}(1 - X^{(0)})$)

The likelihood model and the data – aka. “The information operator”



The likelihood and data relate the prior to the desired posterior: the numerical ODE solution

- **Ideal but intractable goal:** Want $x(t)$ to satisfy the ODE

$$\dot{x}(t) = f(x(t), t)$$

using $X(t)$
 \Leftrightarrow

$$X^{(1)}(t) = f\left(X^{(0)}(t), t\right)$$

$$0 = X^{(1)}(t) - f\left(X^{(0)}(t), t\right) =: m(X(t), t).$$

- **Easier goal:** Satisfy the ODE on a discrete time grid $\{t_i\}_{i=1}^N$

$$\dot{x}(t_i) = f(x(t_i), t_i), \quad i = 1, \dots, N.$$

$$\Leftrightarrow m(X(t_i), t_i) = 0$$

- This motivates the *noiseless likelihood model* and **data**:

$$Z(t_i) | X(t_i) \sim \delta(m(X(t_i), t_i))$$

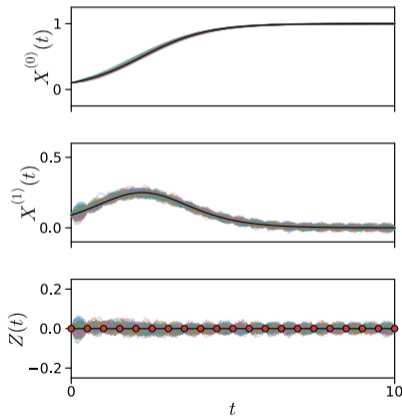
$$z_i \triangleq 0, \quad i = 1, \dots, N.$$

where z_i is a realization of $Z(t_i)$.

(δ is the Dirac distribution)

Example: Logistic ODE $\dot{x} = x(1 - x)$

Posterior samples & ODE solution



(here: $Z = X^{(1)} - X^{(0)}(1 - X^{(0)})$)

Spoiler: This is the thing we want!



How to treat ODEs as the state estimation problem that they really are

$$p \left(x(t) \mid x(0) = x_0, \{ \dot{x}(t_n) = f(x(t_n), t_n) \}_{n=1}^N \right)$$

We want *fast* (approximate) inference \Rightarrow Gaussian filtering and smoothing (it's $\mathcal{O}(N)$!)

1. Prior: q -times integrated Wiener process prior:

$$X(t+h) \mid X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$$

2. **Likelihood:** $Z(t) \mid X(t) \sim \delta(X^{(1)}(t) - f(X^{(0)}(t), t))$

3. **Data:** $\mathcal{D}_{\text{PN}} = \{z_i\}$, with $(Z(t_i) =)z_i = 0$ on a discrete time grid $t_i \in \mathbb{T}$.



Probabilistic numerical ODE solutions

How to treat ODEs as the state estimation problem that they really are

$$p \left(x(t) \mid x(0) = x_0, \{ \dot{x}(t_n) = f(x(t_n), t_n) \}_{n=1}^N \right)$$

We want *fast* (approximate) inference \Rightarrow Gaussian filtering and smoothing (it's $\mathcal{O}(N)$!)

1. Prior: q -times integrated Wiener process prior:

$$X(t+h) \mid X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$$

2. Likelihood: $Z(t) \mid X(t) \sim \delta(X^{(1)}(t) - f(X^{(0)}(t), t))$
3. Data: $\mathcal{D}_{\text{PN}} = \{z_i\}$, with $(Z(t_i) =)z_i = 0$ on a discrete time grid $t_i \in \mathbb{T}$.

This describes a state-space model



Probabilistic numerical ODE solutions

How to treat ODEs as the state estimation problem that they really are

$$p \left(x(t) \mid x(0) = x_0, \{ \dot{x}(t_n) = f(x(t_n), t_n) \}_{n=1}^N \right)$$

We want *fast* (approximate) inference \Rightarrow Gaussian filtering and smoothing (it's $\mathcal{O}(N)$!)

1. Prior: q -times integrated Wiener process prior:

$$X(t+h) \mid X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$$

2. Likelihood: $Z(t) \mid X(t) \sim \delta(X^{(1)}(t) - f(X^{(0)}(t), t))$

3. Data: $\mathcal{D}_{\text{PN}} = \{z_i\}$, with $(Z(t_i) =)z_i = 0$ on a discrete time grid $t_i \in \mathbb{T}$.

This describes a state-space model \Rightarrow solve with EKF/EKS!



The extended Kalman ODE filter – the state-space model

Bringing the last slides all together

For a given initial value problem $\dot{x}(t) = f(x(t), t)$ on $[0, T]$ with $x(0) = x_0$, we have:

The extended Kalman ODE filter – the state-space model

Bringing the last slides all together

For a given initial value problem $\dot{x}(t) = f(x(t), t)$ on $[0, T]$ with $x(0) = x_0$, we have:

Initial distribution:	$X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$
Prior / dynamics model:	$X(t+h) X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$
Likelihood / information model:	$Z(t_i) X(t_i) \sim \delta(X^{(1)}(t_i) - f(X^{(0)}(t_i), t_i))$
Data:	$z_i \triangleq 0, \quad i = 1, \dots, N.$

The extended Kalman ODE filter – the state-space model

Bringing the last slides all together

For a given initial value problem $\dot{x}(t) = f(x(t), t)$ on $[0, T]$ with $x(0) = x_0$, we have:

Initial distribution:	$X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$
Prior / dynamics model:	$X(t+h) X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$
Likelihood / information model:	$Z(t_i) X(t_i) \sim \delta(X^{(1)}(t_i) - f(X^{(0)}(t_i), t_i))$
Data:	$z_i \triangleq 0, \quad i = 1, \dots, N.$

One thing is still missing:

The extended Kalman ODE filter – the state-space model

Bringing the last slides all together

For a given initial value problem $\dot{x}(t) = f(x(t), t)$ on $[0, T]$ with $x(0) = x_0$, we have:

Initial distribution:	$X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$
Prior / dynamics model:	$X(t+h) X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$
Likelihood / information model:	$Z(t_i) X(t_i) \sim \delta(X^{(1)}(t_i) - f(X^{(0)}(t_i), t_i))$
Data:	$z_i \triangleq 0, \quad i = 1, \dots, N.$

One thing is still missing:

What about the initial value??

The extended Kalman ODE filter – the state-space model

Bringing the last slides all together

For a given initial value problem $\dot{x}(t) = f(x(t), t)$ on $[0, T]$ with $x(0) = x_0$, we have:

Initial distribution:	$X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$
Prior / dynamics model:	$X(t+h) X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$
Likelihood / information model:	$Z(t_i) X(t_i) \sim \delta(X^{(1)}(t_i) - f(X^{(0)}(t_i), t_i))$
Data:	$z_i \triangleq 0, \quad i = 1, \dots, N.$

One thing is still missing:

What about the initial value??

Just add another “measurement” at $t = 0$:

$$z^{\text{init}} | X(0) \sim \delta(X^{(0)}(0)), \quad z^{\text{init}} \triangleq x_0.$$

The extended Kalman ODE filter – building blocks

The extended Kalman filter needs these common subroutines

Algorithm 1 Kalman filter prediction

```
1 procedure KF_PREDICT( $\mu, \Sigma, A, Q$ )
2    $\mu^P \leftarrow A\mu$  // Predict mean
3    $\Sigma^P \leftarrow A\Sigma A^T + Q$  // Predict covariance
4   return  $\mu^P, \Sigma^P$ 
5 end procedure
```

Algorithm 2 Extended Kalman filter update

```
1 procedure EKF_UPDATE( $\mu, \Sigma, h, R, y$ )
2    $\hat{y} \leftarrow h(\mu)$  // evaluate the observation model
3    $H \leftarrow J_h(\mu)$  // Jacobian of the observation model
4    $S \leftarrow H\Sigma H^T + R$  // Measurement covariance
5    $K \leftarrow \Sigma H^T S^{-1}$  // Kalman gain
6    $\mu^F \leftarrow \mu + K(y - \hat{y})$  // update mean
7    $\Sigma^F \leftarrow \Sigma - KSK^T$  // update covariance
8   return  $\mu^F, \Sigma^F$ 
9 end procedure
```

(KF_UPDATE analog but with affine h)

The extended Kalman ODE filter

We can solve ODEs with basically just an extended Kalman filter

Algorithm 3 The extended Kalman ODE filter

```

1  procedure EXTENDED KALMAN ODE FILTER( $(\mu_0^-, \Sigma_0^-)$ ,  $(A, Q)$ ,  $(f, x_0)$ ,  $\{t_i\}_{i=1}^N$ )
2       $\mu_0, \Sigma_0 \leftarrow$  KF_UPDATE( $\mu_0^-, \Sigma_0^-, E_0, \mathbf{0}_{d \times d}, x_0$ )           // Initial update to fit the initial value
3      for  $k \in \{1, \dots, N\}$  do
4           $h_k \leftarrow t_k - t_{k-1}$                                            // Step size
5           $\mu_k^-, \Sigma_k^- \leftarrow$  KF_PREDICT( $\mu_{k-1}, \Sigma_{k-1}, A(h_k), Q(h_k)$ ) // Kalman filter prediction
6           $m_k(X) := E_1 X - f(E_0 X, t_k)$                                      // Define the non-linear observation model
7           $\mu_k, \Sigma_k \leftarrow$  EKF_UPDATE( $\mu_k^-, \Sigma_k^-, m_k, \mathbf{0}_{d \times d}, \mathbf{0}_d$ ) // Extended Kalman filter update
8      end for
9      return  $(\mu_k, \Sigma_k)_{k=1}^N$ 
10 end procedure
    
```

Recall: The *state* $X(t)$ is a stack of q derivatives $X = [X^{(0)}, X^{(1)}, \dots, X^{(q)}]^T$.

The projection matrices E_i map X to the i -th derivative: $E_i X = X^{(i)}$.

The extended Kalman ODE filter

We can solve ODEs with basically just an extended Kalman filter

Algorithm 3 The extended Kalman ODE filter

```

1  procedure EXTENDED KALMAN ODE FILTER( $(\mu_0^-, \Sigma_0^-)$ ,  $(A, Q)$ ,  $(f, x_0)$ ,  $\{t_i\}_{i=1}^N$ )
2       $\mu_0, \Sigma_0 \leftarrow$  KF_UPDATE( $\mu_0^-, \Sigma_0^-, E_0, \mathbf{0}_{d \times d}, x_0$ ) // Initial update to fit the initial value
3      for  $k \in \{1, \dots, N\}$  do
4           $h_k \leftarrow t_k - t_{k-1}$  // Step size
5           $\mu_k^-, \Sigma_k^- \leftarrow$  KF_PREDICT( $\mu_{k-1}, \Sigma_{k-1}, A(h_k), Q(h_k)$ ) // Kalman filter prediction
6           $m_k(X) := E_1 X - f(E_0 X, t_k)$  // Define the non-linear observation model
7           $\mu_k, \Sigma_k \leftarrow$  EKF_UPDATE( $\mu_k^-, \Sigma_k^-, m_k, \mathbf{0}_{d \times d}, \mathbf{0}_d$ ) // Extended Kalman filter update
8      end for
9      return  $(\mu_k, \Sigma_k)_{k=1}^N$ 
10 end procedure
    
```

Recall: The *state* $X(t)$ is a stack of q derivatives $X = [X^{(0)}, X^{(1)}, \dots, X^{(q)}]^T$.
 The projection matrices E_i map X to the i -th derivative: $E_i X = X^{(i)}$.

EXTENDED KALMAN ODE SMOOTHER: Just run a RTS smoother after the filter!

DEMO TIME: The extended Kalman ODE filter in code

`demo.jl`



Uncertainty calibration or “how to choose prior hyperparameters”

Hyperparameters of the prior have a strong influence on posteriors – so we need to estimate them

- ▶ Recall the IWP(q) prior model: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), \sigma^2 Q(h))$.
⇒ The hyperparameter σ directly influences covariances! But what value should it have?

Uncertainty calibration or “how to choose prior hyperparameters”

Hyperparameters of the prior have a strong influence on posteriors – so we need to estimate them

- ▶ Recall the IWP(q) prior model: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), \sigma^2 Q(h))$.
⇒ The hyperparameter σ directly influences covariances! But what value should it have?
- ▶ **Standard approach:** Maximize the marginal likelihood:

$$\hat{\sigma} = \arg \max p(\mathcal{D}_{\text{PN}} | \sigma) = p(z_{1:N} | \sigma) = p(z_1 | \sigma) \prod_{k=2}^N p(z_k | z_{1:k-1}, \sigma).$$

Uncertainty calibration or “how to choose prior hyperparameters”

Hyperparameters of the prior have a strong influence on posteriors – so we need to estimate them

- ▶ Recall the IWP(q) prior model: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), \sigma^2 Q(h))$.
⇒ The hyperparameter σ directly influences covariances! But what value should it have?
- ▶ **Standard approach:** Maximize the marginal likelihood:

$$\hat{\sigma} = \arg \max p(\mathcal{D}_{\text{PN}} | \sigma) = p(z_{1:N} | \sigma) = p(z_1 | \sigma) \prod_{k=2}^N p(z_k | z_{1:k-1}, \sigma).$$

- ▶ The EKF provides Gaussian estimates $p(z_k | z_{1:k-1}) \approx \mathcal{N}(z_k; \hat{z}_k, S_k)$. This gives a *quasi*-MLE:

$$\hat{\sigma} = \arg \max p(\mathcal{D}_{\text{PN}} | \sigma) \approx \arg \max \sum_{k=1}^N \log \mathcal{N}(z_k; \hat{z}_k, S_k).$$

Uncertainty calibration or “how to choose prior hyperparameters”

Hyperparameters of the prior have a strong influence on posteriors – so we need to estimate them

- ▶ Recall the IWP(q) prior model: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), \sigma^2 Q(h))$.
 \Rightarrow The hyperparameter σ directly influences covariances! But what value should it have?
- ▶ **Standard approach:** Maximize the marginal likelihood:

$$\hat{\sigma} = \arg \max p(\mathcal{D}_{\text{PN}} | \sigma) = p(z_{1:N} | \sigma) = p(z_1 | \sigma) \prod_{k=2}^N p(z_k | z_{1:k-1}, \sigma).$$

- ▶ The EKF provides Gaussian estimates $p(z_k | z_{1:k-1}) \approx \mathcal{N}(z_k; \hat{z}_k, S_k)$. This gives a *quasi*-MLE:

$$\hat{\sigma} = \arg \max p(\mathcal{D}_{\text{PN}} | \sigma) \approx \arg \max \sum_{k=1}^N \log \mathcal{N}(z_k; \hat{z}_k, S_k).$$

- ▶ **In our specific context this can be solved in closed form:**

$$\hat{\sigma}^2 = \frac{1}{Nd} \sum_{i=1}^N (z_i - \hat{z}_i)^\top S_i^{-1} (z_i - \hat{z}_i).$$

We don't even need to run the filter again! Just adjust the covariances: (proof: homework)

$$\Sigma_i \leftarrow \hat{\sigma}^2 \cdot \Sigma_i, \quad \forall i \in \{1, \dots, N\}.$$

DEMO TIME: Calibrated vs uncalibrated posteriors

demo.jl





Why?



Why be probabilistic about ODE solutions?



Why be probabilistic about ODE solutions?

- ▶ **Uncertainty quantification:**

The methods provide estimates of their numerical error



Why be probabilistic about ODE solutions?

- ▶ **Uncertainty quantification:**

The methods provide estimates of their numerical error

- ▶ **Flexibility / convenience / efficiency:**

The probabilistic state-space formulation makes it very easy to perform joint inference on various kinds of information (this is what we will do next!)



Example 1: Extending ODE filters to other related problems by adjusting the *information model*



Numerical problem setting: Initial value problem with ODE

$$\dot{x}(t) = f(x(t), t), \quad x(0) = x_0.$$

This leads to the **probabilistic state estimation problem:**

Initial distribution: $X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$

Prior / dynamics model: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$

ODE likelihood: $Z(t_i) | X(t_i) \sim \delta\left(X^{(1)}(t_i) - f(X^{(0)}(t_i), t_i)\right), \quad z_i \triangleq 0$

Initial value likelihood: $Z^{\text{init}} | X(0) \sim \delta\left(X^{(0)}(0)\right), \quad z^{\text{init}} \triangleq x_0$

Numerical problem setting: Initial value problem with **second-order** ODE

$$\ddot{x}(t) = f(\dot{x}(t), x(t), t), \quad x(0) = x_0, \quad \dot{x}(0) = \dot{x}_0.$$

This leads to the **probabilistic state estimation problem**:

Initial distribution: $X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$

Prior / dynamics model: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$

ODE likelihood: $Z(t_i) | X(t_i) \sim \delta\left(X^{(1)}(t_i) - f(X^{(0)}(t_i), t_i)\right), \quad z_i \triangleq 0$

Initial value likelihood: $Z^{\text{init}} | X(0) \sim \delta\left(X^{(0)}(0)\right), \quad z^{\text{init}} \triangleq x_0$

Numerical problem setting: Initial value problem with **second-order** ODE

$$\ddot{x}(t) = f(\dot{x}(t), x(t), t), \quad x(0) = x_0, \quad \dot{x}(0) = \dot{x}_0.$$

This leads to the **probabilistic state estimation problem**:

Initial distribution: $X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$

Prior / dynamics model: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$

ODE likelihood: $Z(t_i) | X(t_i) \sim \delta \left(X^{(2)}(t_i) - f(X^{(1)}(t_i), X^{(0)}(t_i), t_i) \right), \quad z_i \triangleq 0$

Initial value likelihood: $Z^{\text{init}} | X(0) \sim \delta \left(X^{(0)}(0) \right), \quad z^{\text{init}} \triangleq x_0$

Initial derivative likelihood: $Z_1^{\text{init}} | X(0) \sim \delta \left(X^{(1)}(0) \right), \quad z_1^{\text{init}} \triangleq \dot{x}_0$

Numerical problem setting: Initial value problem with *differential-algebraic equation* (DAE) in mass-matrix form

$$M\dot{x}(t) = f(x(t), t), \quad x(0) = x_0. \quad (\text{with singular } M)$$

This leads to the **probabilistic state estimation problem**:

Initial distribution: $X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$

Prior / dynamics model: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$

ODE likelihood: $Z(t_j) | X(t_j) \sim \delta(X^{(1)}(t_j) - f(X^{(0)}(t_j), t_j)), \quad z_j \triangleq 0$

Initial value likelihood: $z^{\text{init}} | X(0) \sim \delta(X^{(0)}(0)), \quad z^{\text{init}} \triangleq x_0$

Numerical problem setting: Initial value problem with *differential-algebraic equation* (DAE) in mass-matrix form

$$M\dot{x}(t) = f(x(t), t), \quad x(0) = x_0. \quad (\text{with singular } M)$$

This leads to the **probabilistic state estimation problem**:

Initial distribution: $X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$

Prior / dynamics model: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$

DAE likelihood: $Z(t_i) | X(t_i) \sim \delta\left(MX^{(1)}(t_i) - f(X^{(0)}(t_i), t_i)\right), \quad z_i \triangleq 0$

Initial value likelihood: $Z^{\text{init}} | X(0) \sim \delta\left(X^{(0)}(0)\right), \quad z^{\text{init}} \triangleq x_0$

Numerical problem setting: Initial value problem with first-order ODE and **conserved quantities**

$$\dot{x}(t) = f(x(t), t), \quad x(0) = x_0, \quad g(x(t), \dot{x}(t)) = 0.$$

This leads to the **probabilistic state estimation problem**:

Initial distribution: $X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$

Prior / dynamics model: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$

ODE likelihood: $Z(t_i) | X(t_i) \sim \delta\left(X^{(1)}(t_i) - f(X^{(0)}(t_i), t_i)\right), \quad z_i \triangleq 0$

Initial value likelihood: $Z^{\text{init}} | X(0) \sim \delta\left(X^{(0)}(0)\right), \quad z^{\text{init}} \triangleq x_0$

Numerical problem setting: Initial value problem with first-order ODE and **conserved quantities**

$$\dot{x}(t) = f(x(t), t), \quad x(0) = x_0, \quad g(x(t), \dot{x}(t)) = 0.$$

This leads to the **probabilistic state estimation problem**:

Initial distribution: $X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$

Prior / dynamics model: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$

ODE likelihood: $Z(t_i) | X(t_i) \sim \delta\left(X^{(1)}(t_i) - f(X^{(0)}(t_i), t_i)\right), \quad z_i \triangleq 0$

Conservation law likelihood: $Z_i^c(t_i) | X(t_i) \sim \delta\left(g(X^{(0)}(t_i), X^{(1)}(t_i))\right), \quad z_i^c \triangleq 0$

Initial value likelihood: $Z^{\text{init}} | X(0) \sim \delta\left(X^{(0)}(0)\right), \quad z^{\text{init}} \triangleq x_0$



DEMO TIME: Solving a second-order ODE

demo.jl



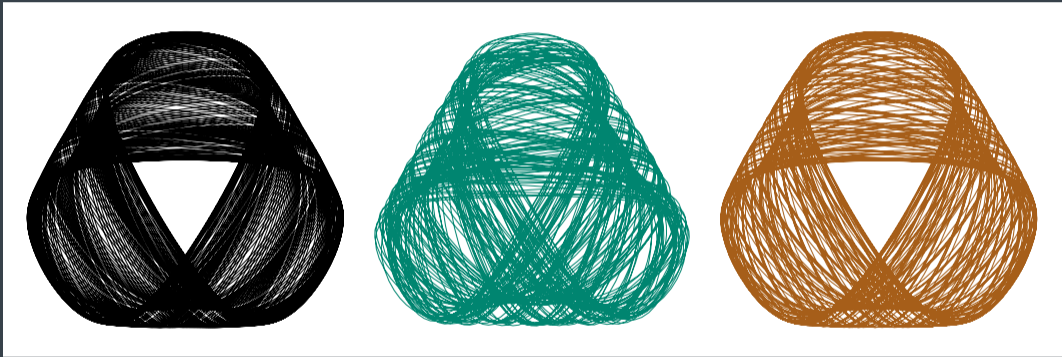
DEMO TIME: Conserved quantities

henonheiles.mp4





DEMO TIME: Conserved quantities





Example 2: Combine ODEs and GP regression via *latent force inference*

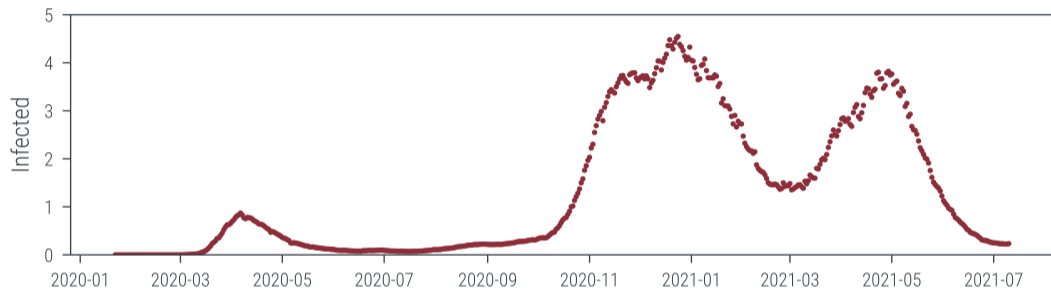


Latent force inference: GP regression on both ODEs and data



An example we know all too well: COVID-19

Paper: *Schmidt, Krämer, Hennig, NeurIPS 2021*



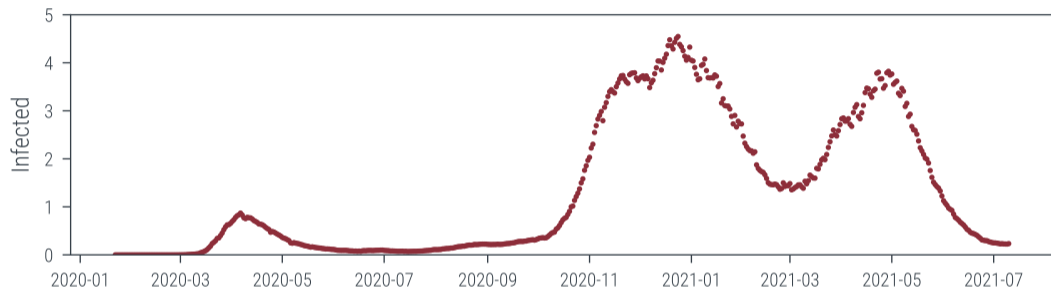
ODE dynamics:

$$\frac{d}{dt}x(t) = f(x(t), t)$$

Latent force inference: GP regression on both ODEs and data

An example we know all too well: COVID-19

Paper: *Schmidt, Krämer, Hennig, NeurIPS 2021*



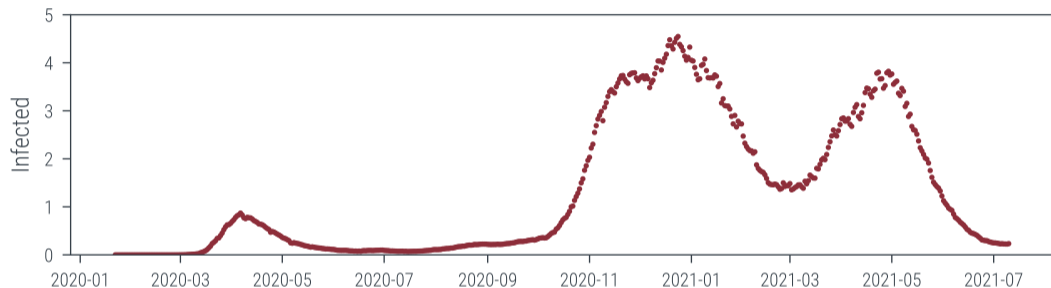
ODE dynamics:

$$\frac{d}{dt} \begin{bmatrix} S(t) \\ I(t) \\ R(t) \\ D(t) \end{bmatrix} = \begin{bmatrix} -\beta \cdot S(t)I(t)/P \\ \beta \cdot S(t)I(t)/P - \gamma I(t) - \eta I(t) \\ \gamma I(t) \\ \eta I(t) \end{bmatrix}$$

Latent force inference: GP regression on both ODEs and data

An example we know all too well: COVID-19

Paper: *Schmidt, Krämer, Hennig, NeurIPS 2021*



ODE dynamics with time-varying contact rate:

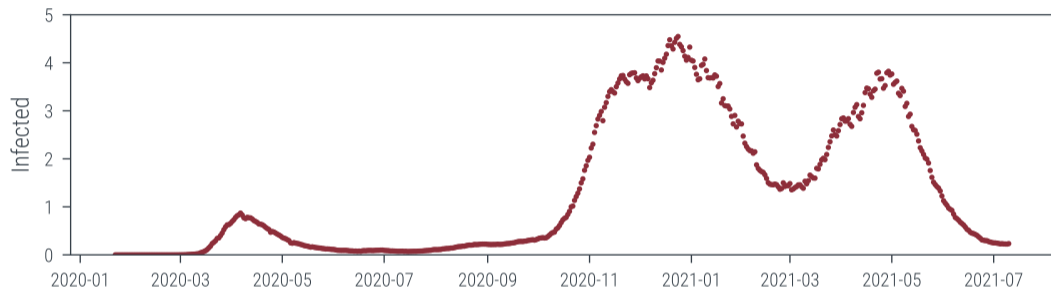
$$\frac{d}{dt} \begin{bmatrix} S(t) \\ I(t) \\ R(t) \\ D(t) \end{bmatrix} = \begin{bmatrix} -\beta(t) \cdot S(t)I(t)/P \\ \beta(t) \cdot S(t)I(t)/P - \gamma I(t) - \eta I(t) \\ \gamma I(t) \\ \eta I(t) \end{bmatrix}$$

Latent force inference: GP regression on both ODEs and data



An example we know all too well: COVID-19

Paper: *Schmidt, Krämer, Hennig, NeurIPS 2021*



ODE dynamics with time-varying contact rate:

$$\frac{d}{dt} \begin{bmatrix} S(t) \\ I(t) \\ R(t) \\ D(t) \end{bmatrix} = \begin{bmatrix} -\beta(t) \cdot S(t)I(t)/P \\ \beta(t) \cdot S(t)I(t)/P - \gamma I(t) - \eta I(t) \\ \gamma I(t) \\ \eta I(t) \end{bmatrix}$$

Latent force model: Gauss–Markov process

$$\beta(t+h) \mid \beta(t) \sim \mathcal{N}(A_\beta(h)\beta(t), Q_\beta(h))$$

Data:

$$y_i \mid x(t_i) \sim \mathcal{N}(Hx(t_i), \sigma^2 I)$$

Formally we obtain the **probabilistic state estimation problem**:

State initial distribution: $X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$

State dynamics: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$

Latent force initial distribution: $\beta(0) \sim \mathcal{N}(\mu_0^\beta, \Sigma_0^\beta)$

Latent force dynamics: $\beta(t+h) | \beta(t) \sim \mathcal{N}(A_\beta(h)\beta(t), Q_\beta(h))$

ODE likelihood: $Z(t_i) | X(t_i), \beta(t_i) \sim \delta(X^{(1)}(t_i) - f(X^{(0)}(t_i), \beta(t_i), t_i)), \quad z_i \triangleq 0$

Initial value likelihood: $Z^{\text{init}} | X(0) \sim \delta(X^{(0)}(0)), \quad z^{\text{init}} \triangleq x_0$

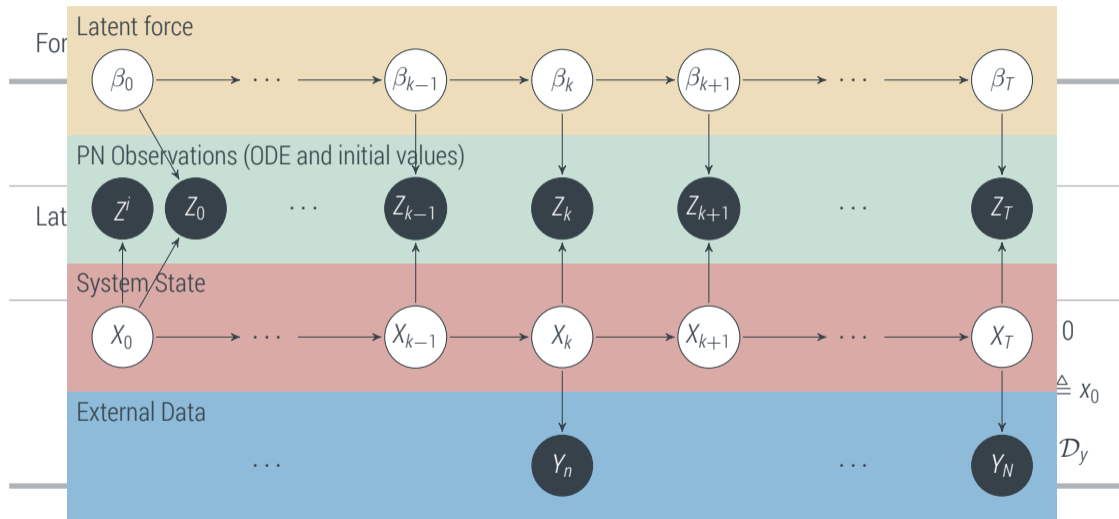
Data likelihood: $Y_i | X(t_i) \sim \mathcal{N}(HX^{(0)}(t_i), \sigma^2 I), \quad y_i \in \mathcal{D}_y$

Latent force inference: Writing down the state estimation problem



Once again it's just a state estimation problem

Paper: Schmidt, Krämer, Hennig, NeurIPS 2021



Formally we obtain the **probabilistic state estimation problem**:

State initial distribution: $X(0) \sim \mathcal{N}(\mu_0, \Sigma_0)$

State dynamics: $X(t+h) | X(t) \sim \mathcal{N}(A(h)X(t), Q(h))$

Latent force initial distribution: $\beta(0) \sim \mathcal{N}(\mu_0^\beta, \Sigma_0^\beta)$

Latent force dynamics: $\beta(t+h) | \beta(t) \sim \mathcal{N}(A_\beta(h)\beta(t), Q_\beta(h))$

ODE likelihood: $Z(t_i) | X(t_i), \beta(t_i) \sim \delta(X^{(1)}(t_i) - f(X^{(0)}(t_i), \beta(t_i), t_i)), \quad z_i \triangleq 0$

Initial value likelihood: $Z^{\text{init}} | X(0) \sim \delta(X^{(0)}(0)), \quad z^{\text{init}} \triangleq x_0$

Data likelihood: $Y_i | X(t_i) \sim \mathcal{N}(HX^{(0)}(t_i), \sigma^2 I), \quad y_i \in \mathcal{D}_y$

Latent force inference: Writing down the state estimation problem

Once again it's just a state estimation problem

Paper: *Schmidt, Krämer, Hennig, NeurIPS 2021*

Simplify by stacking $\tilde{X} = [X, \beta]$:

Initial distribution: $\tilde{X}(0) \sim \mathcal{N}(\tilde{\mu}_0, \tilde{\Sigma}_0)$

Prior / dynamics model: $\tilde{X}(t+h) | \tilde{X}(t) \sim \mathcal{N}(\tilde{A}(h)\tilde{X}(t), \tilde{Q}(h))$

ODE likelihood: $Z(t_i) | \tilde{X}(t_i) \sim \delta(E_1\tilde{X}(t_i) - f(E_0\tilde{X}(t_i), E_\beta\tilde{X}(t_i), t_i)), \quad z_i \triangleq 0$

Initial value likelihood: $Z^{\text{init}} | \tilde{X}(0) \sim \delta(E_0\tilde{X}(0)), \quad z^{\text{init}} \triangleq x_0$

Data likelihood: $Y_i | \tilde{X}(t_i) \sim \mathcal{N}(HE_0\tilde{X}(t_i), \sigma^2 I), \quad y_i \in \mathcal{D}_y$

with $E_0\tilde{X} := X^{(0)}, E_1\tilde{X} := X^{(1)}, E_\beta\tilde{X} := \beta$.

Simplify by stacking $\tilde{X} = [X, \beta]$:

Initial distribution: $\tilde{X}(0) \sim \mathcal{N}(\tilde{\mu}_0, \tilde{\Sigma}_0)$

Prior / dynamics model: $\tilde{X}(t+h) | \tilde{X}(t) \sim \mathcal{N}(\tilde{A}(h)\tilde{X}(t), \tilde{Q}(h))$

ODE likelihood: $Z(t_i) | \tilde{X}(t_i) \sim \delta(E_1\tilde{X}(t_i) - f(E_0\tilde{X}(t_i), E_\beta\tilde{X}(t_i), t_i)), \quad z_i \triangleq 0$

Initial value likelihood: $Z^{\text{init}} | \tilde{X}(0) \sim \delta(E_0\tilde{X}(0)), \quad z^{\text{init}} \triangleq x_0$

Data likelihood: $Y_i | \tilde{X}(t_i) \sim \mathcal{N}(HE_0\tilde{X}(t_i), \sigma^2 I), \quad y_i \in \mathcal{D}_y$

with $E_0\tilde{X} := X^{(0)}, E_1\tilde{X} := X^{(1)}, E_\beta\tilde{X} := \beta$.

Again: **This is just a state-space model**

Latent force inference: Writing down the state estimation problem

Once again it's just a state estimation problem

Paper: *Schmidt, Krämer, Hennig, NeurIPS 2021*

Simplify by stacking $\tilde{X} = [X, \beta]$:

Initial distribution: $\tilde{X}(0) \sim \mathcal{N}(\tilde{\mu}_0, \tilde{\Sigma}_0)$

Prior / dynamics model: $\tilde{X}(t+h) | \tilde{X}(t) \sim \mathcal{N}(\tilde{A}(h)\tilde{X}(t), \tilde{Q}(h))$

ODE likelihood: $Z(t_i) | \tilde{X}(t_i) \sim \delta(E_1\tilde{X}(t_i) - f(E_0\tilde{X}(t_i), E_\beta\tilde{X}(t_i), t_i)), \quad z_i \triangleq 0$

Initial value likelihood: $Z^{\text{init}} | \tilde{X}(0) \sim \delta(E_0\tilde{X}(0)), \quad z^{\text{init}} \triangleq x_0$

Data likelihood: $Y_i | \tilde{X}(t_i) \sim \mathcal{N}(HE_0\tilde{X}(t_i), \sigma^2 I), \quad y_i \in \mathcal{D}_y$

with $E_0\tilde{X} := X^{(0)}, E_1\tilde{X} := X^{(1)}, E_\beta\tilde{X} := \beta$.

Again: **This is just a state-space model \Rightarrow inference with EKF/EKS!**

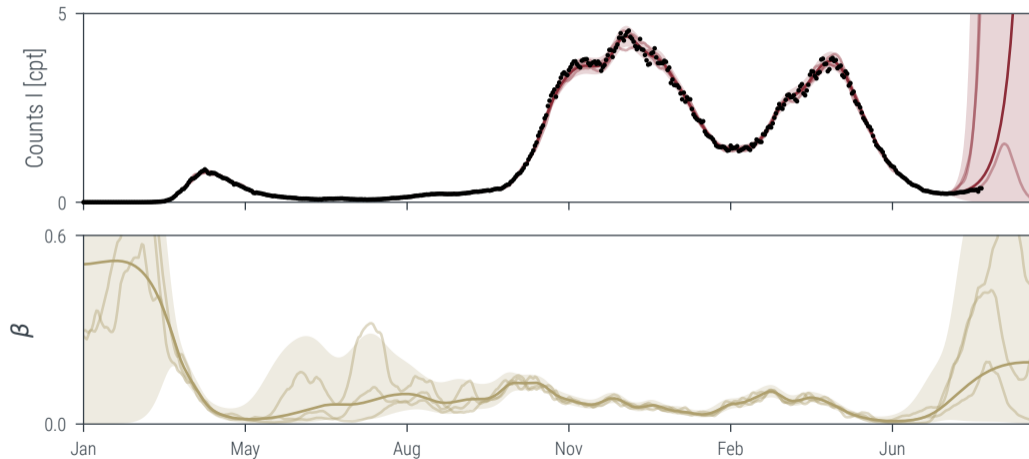
Latent force inference: Results



Posteriors over infections and contact rates *in a single forward-backward pass*

Paper: Schmidt, Krämer, Hennig, NeurIPS 2021

The extended Kalman smoother returns probabilistic estimates for all states (S, I, R, D, β):



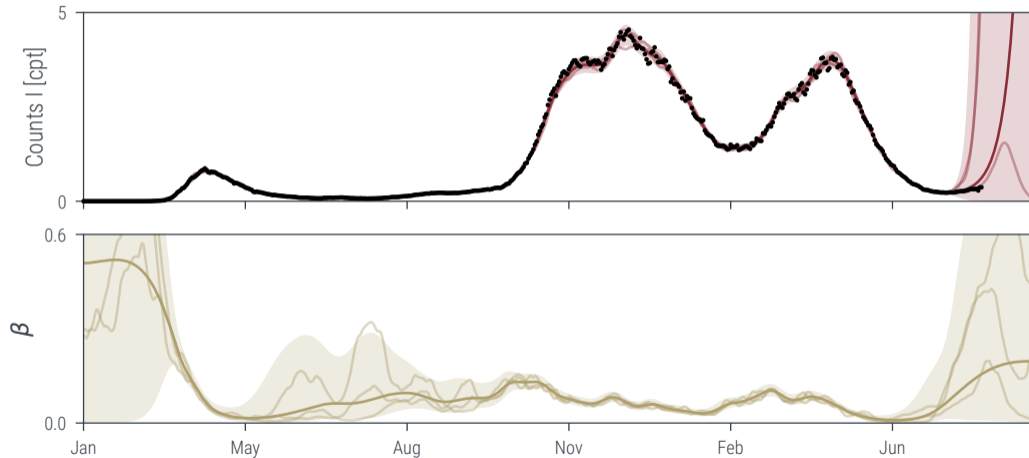
Latent force inference: Results



Posteriors over infections and contact rates *in a single forward-backward pass*

Paper: Schmidt, Krämer, Hennig, NeurIPS 2021

The extended Kalman smoother returns probabilistic estimates for all states (S, I, R, D, β):



⇒ Probabilistic estimates of a latent force *in a single forward-backward pass!*

- ▶ ODE solving is state estimation
⇒ treat ODEs as state estimation problems!
- ▶ **We can solve ODEs with Bayesian filtering and smoothing**
⇒ “ODE filters”
- ▶ *Flexible information operators:*
Easily adjust the model to solve other numerical problems,
with essentially the same algorithm!
- ▶ *Latent force inference:*
Filters enable *efficient* joint inference on both ODEs and data.

Please cite this course, as

```
@techreport{NoML22,  
  title = {Numerics of Machine Learning},  
  author = {N. Bosch and J. Grosse  
and P. Hennig and A. Kristiadi  
and M. Pförtner and J. Schmidt  
and F. Schneider and L. Tatzel  
and J. Wenger},  
  series = {Lecture Notes in Machine Learning},  
  year = {2022},  
  institution = {Tübingen AI Center},  
}
```

Bayesian filtering and smoothing is the right framework for modeling dynamical systems
in a modular and data-centric fashion.

Next week: *Partial* Differential Equations!

