



OST

Ostschweizer
Fachhochschule

Introduction to recursive probabilistic learning

Applications

JuanPi Carbajal
15.10.2021

Institute for Energy Technology

Recap and Q&A

Kalman filter

$$\begin{array}{l} \text{while no measurement: predict} \\ \text{predict measurement} \\ \text{on measurement: update} \end{array} \left\{ \begin{array}{l} \hat{\mathbf{m}}_k = \mathbf{A}_k \mathbf{m}_{k-1} \\ \hat{\mathbf{P}}_k = \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^\top + \mathbf{Q}_{k-1} \\ \hat{\mathbf{y}}_k = \mathbf{H}_k \hat{\mathbf{m}}_k \\ \hat{\mathbf{S}}_k = \mathbf{H}_k \hat{\mathbf{P}}_k \mathbf{H}_k^\top + \mathbf{R}_k \\ \mathbf{K}_k = \hat{\mathbf{P}}_k \mathbf{H}_k^\top \hat{\mathbf{S}}_k^{-1} \\ \mathbf{m}_k = \hat{\mathbf{m}}_k + \mathbf{K}_k (\mathbf{y}_k - \hat{\mathbf{y}}_k) \\ \mathbf{P}_k = \hat{\mathbf{P}}_k - \mathbf{K}_k \hat{\mathbf{S}}_k \mathbf{K}_k^\top \end{array} \right.$$

Do not forget to update your measurement prediction after and update.

Python & GNU Octave (MATLAB) tools

These libraries are relatively easy to follow, and have a book for documentation:

- FilterPy
 - <https://filterpy.readthedocs.io/en/latest/>
 - book: https://nbviewer.org/github/rlabbe/Kalman-and-Bayesian-Filters-in-Python/blob/master/table_of_contents.ipynb
- ekfukf
 - MATLAB/Octave <https://github.com/EEA-sensors/ekfukf>,
 - Octave installable <https://github.com/kakila/ekfukf>
 - book http://users.aalto.fi/~ssarkka/pub/cup_book_online_20131111.pdf (if you like it, buy it!)

They are not performant (pykalman is abandoned), for prototyping. General probabilistic inference: PyMC3 <https://docs.pymc.io/en/stable/>

There are many libraries for Kalman filter, search for the one that suits you.

In my experience, I use simple ones for quickly prototyping, then implement what I need case specific. Algorithms are very simple. Difficulty in inverting matrices: good algorithms are already available. Most relevant matrices are symmetric positive definite.

Consider Julia <https://julialang.org/>

Warm-up

- Implement recursive linear regression $y_k = t_k a + b + \epsilon_k$
- Delayed linear regression $y_k = w_1 y_{k-1} + \dots + w_d y_{k-d} + \epsilon_k$

Discretization of ODEs

What's an ODE?

Ordinary differential equations (ODE) define a function (the unknown) via a relation of its derivatives:

$$\frac{d^2x(t)}{dt^2} + \gamma \frac{dx(t)}{dt} + \omega^2 x(t) = w(t)$$

That's a **2nd order** equation. For low order we can shorten it to:

$$\ddot{x} + \gamma \dot{x} + \omega^2 x = w(t)$$

High order equations can (almost always) be converted to a **system of 1st order** equations:

$$x_1(t) := x(t), \quad x_2(t) := \dot{x}(t)$$
$$\dot{\mathbf{x}} = \begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -\omega^2 & -\gamma \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} w(t) = \mathbf{F}\mathbf{x} + \mathbf{L}w(t)$$

We exchange order for dimension.

Linear ODE

In general \mathbf{F} and \mathbf{L} can depend on \mathbf{x} and time. If they are constant or depend only on time the system is **linear** (affine):

$$\dot{\mathbf{x}} = \mathbf{F}(t)\mathbf{x} + \mathbf{L}(t)w_x(t)$$

$$\dot{\mathbf{y}} = \mathbf{F}(t)\mathbf{y} + \mathbf{L}(t)w_y(t)$$

$$\mathbf{z} := \mathbf{x} + a\mathbf{y}$$

$$\begin{aligned}\dot{\mathbf{z}} &= \dot{\mathbf{x}} + a\dot{\mathbf{y}} = \mathbf{F}(t)\mathbf{x} + \mathbf{L}(t)w_x(t) + a\mathbf{F}(t)\mathbf{y} + a\mathbf{L}(t)w_y(t) = \\ &= \mathbf{F}(t)(\mathbf{x} + a\mathbf{y}) + \mathbf{L}(t)(w_x(t) + aw_y(t))\end{aligned}$$

The derivative of the combination follows the same equation with the combined inputs.

1st order homogenous ODE

Defines a 1-dimensional system: $\dot{x} = Fx$. Integrating both sides from 0 to t .

$$\int_0^t \dot{x}(\tau) d\tau = x(t) - x(0) \rightarrow x(t) = x(0) + \int_0^t Fx(\tau) d\tau$$

Recursively insert the formal solution in its expression:

$$\begin{aligned} x(t) &= x(0) + \int_0^t F [x(0) + \int_0^\tau Fx(\tau_1) d\tau_1] d\tau = \\ &= x(0) + Fx(0)t + \int_0^t \int_0^\tau F^2 x(\tau_1) d\tau_1 d\tau = \\ &= x(0) + Fx(0)t + \int_0^t \int_0^\tau F^2 [x(0) + \int_0^{\tau_1} Fx(\tau_2) d\tau_2] d\tau_1 d\tau = \\ &= x(0) + Fx(0)t + \frac{1}{2} F^2 x(0)t^2 + \int_0^t \int_0^\tau \int_0^{\tau_1} F^3 x(\tau_2) d\tau_2 d\tau_1 d\tau = \dots \\ x(t) &= x(0) + Fx(0)t + \frac{1}{2} F^2 x(0)t^2 + \frac{1}{6} F^3 x(0)t^3 + \dots = \\ &= \left(1 + Ft + \frac{F^2 t^2}{2!} + \frac{F^3 t^3}{3!} + \dots \right) x(0) = \\ &= \exp(Ft)x(0) \end{aligned}$$

The infinite sum inside the parenthesis is the Taylor expansion of $\exp(Ft)$. Check it by expanding e^x around 0, i.e.

$$f(x) = \sum_{n=0}^{\infty} \frac{1}{n!} \left. \frac{d^n f}{dx^n} \right|_{x_0} (x - x_0)^n$$

is the Taylor expansion of $f(x)$ around x_0

1st order N-dimensional ODE

The homogenous case, defines a N-dimensional system: $\dot{\mathbf{x}} = \mathbf{F}\mathbf{x}$. Following the same procedure as before

$$\mathbf{x}(t) = \left(\mathbf{I} + \mathbf{F}t + \frac{\mathbf{F}^2 t^2}{2!} + \frac{\mathbf{F}^3 t^3}{3!} + \dots \right) \mathbf{x}(0) := \exp(\mathbf{F}t)\mathbf{x}(0)$$

That's the definition of the **matrix exponential**

The matrix exponential is not the exponential of each entry in the matrix (element-wise exponential). Consider

$$\begin{aligned} \mathbf{F} &= \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \\ \mathbf{F}^2 &= \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} = \mathbf{0} \\ \rightarrow \mathbf{F}^n &= \mathbf{0} \end{aligned}$$

Then

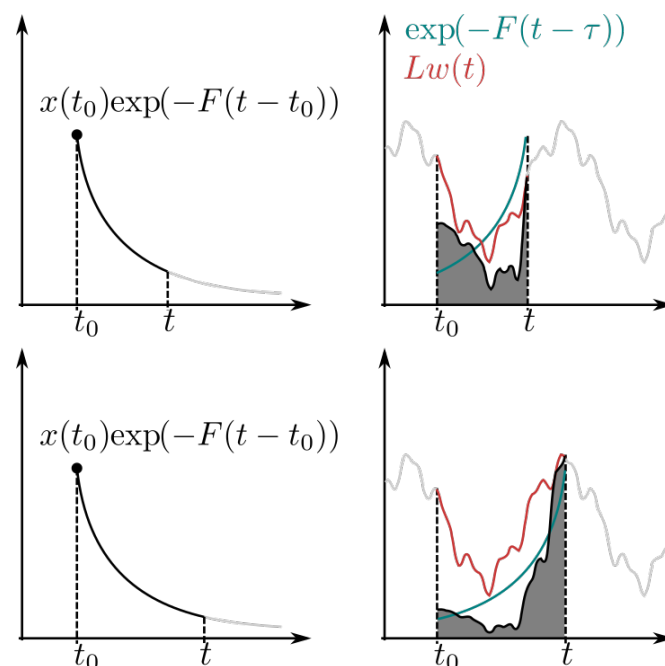
$$\exp(\mathbf{F}t) = \mathbf{I} + \mathbf{F}t = \begin{bmatrix} 1 & t \\ 0 & 1 \end{bmatrix} \neq \begin{bmatrix} e^0 & e^t \\ e^0 & e^0 \end{bmatrix} = \begin{bmatrix} 1 & e^t \\ 1 & 1 \end{bmatrix}$$

1st order N-dimensional ODE

For $\dot{\mathbf{x}} - \mathbf{F}\mathbf{x} = \mathbf{L}\mathbf{w}(t)$. We get the solution between t_0 and t :

$$\mathbf{x}(t) = \exp(\mathbf{F}(t - t_0))\mathbf{x}(0) + \int_{t_0}^t \exp(\mathbf{F}(t - \tau))\mathbf{L}\mathbf{w}(\tau)d\tau$$

In the scalar case this looks like



Multiply both sides by $\exp(-\mathbf{F}t)$

$$\exp(-\mathbf{F}t)\dot{\mathbf{x}} - \exp(-\mathbf{F}t)\mathbf{F}\mathbf{x} = \exp(-\mathbf{F}t)\mathbf{L}\mathbf{w}(t)$$

You can prove from the definition of the matrix exponential the property

$$\frac{d}{dt} \exp(-\mathbf{F}t) = -\exp(-\mathbf{F}t)\mathbf{F}$$

Just like the scalar case, but here the order of multiplication is important. Observe that:

$$\frac{d}{dt} [\exp(-\mathbf{F}t)\mathbf{x}(t)] = \exp(-\mathbf{F}t)\dot{\mathbf{x}} - \exp(-\mathbf{F}t)\mathbf{F}(t)\mathbf{x}$$

Then the equation reduces to

$$\frac{d}{dt} [\exp(-\mathbf{F}t)\mathbf{x}(t)] = \exp(-\mathbf{F}t)\mathbf{L}(t)\mathbf{w}(t)$$

and by integrating both sides from t_0 to t we get the solution. Note that

$$\int_{t_0}^t \frac{d}{d\tau} [\exp(-\mathbf{F}\tau)\mathbf{x}(\tau)]d\tau = \exp(-\mathbf{F}t)\mathbf{x}(t) - \exp(-\mathbf{F}t_0)\mathbf{x}(t_0)$$

Discretization

For the general ODE

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

If we move from t to $t + \Delta t$ we get the solution

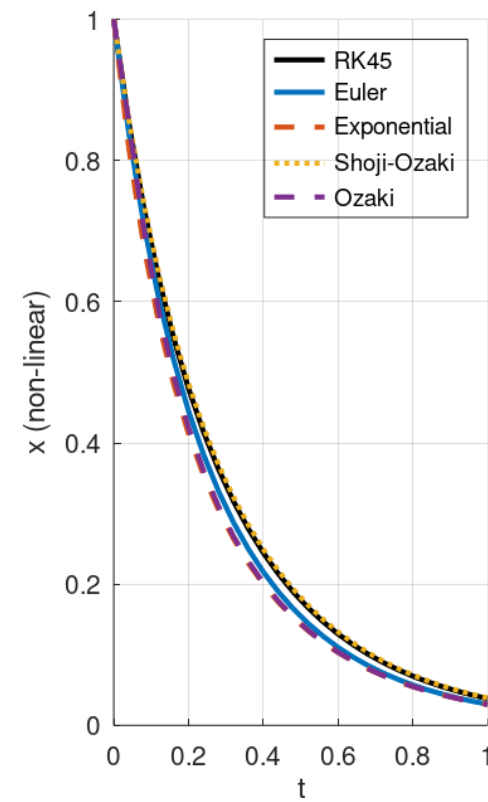
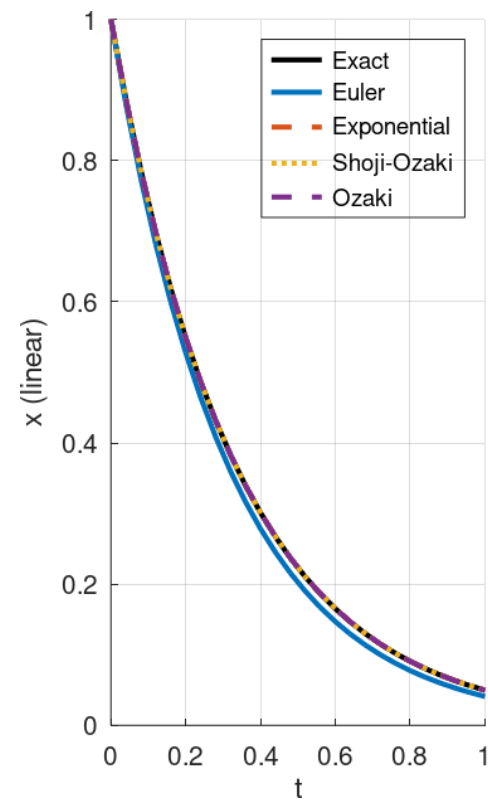
$$\mathbf{x}(t + \Delta t) = \mathbf{x}(t) + \int_t^{t+\Delta t} \mathbf{f}(\mathbf{x}(\tau), \tau) d\tau$$

So the question is how to discretize the integral in the right-hand side. Choosing a quadrature defines an integration method. Euler method is obtained by the approximation

$$\int_t^{t+\Delta t} \mathbf{f}(\mathbf{x}(\tau), \tau) d\tau \approx \mathbf{f}(\mathbf{x}(t), t) \Delta t$$

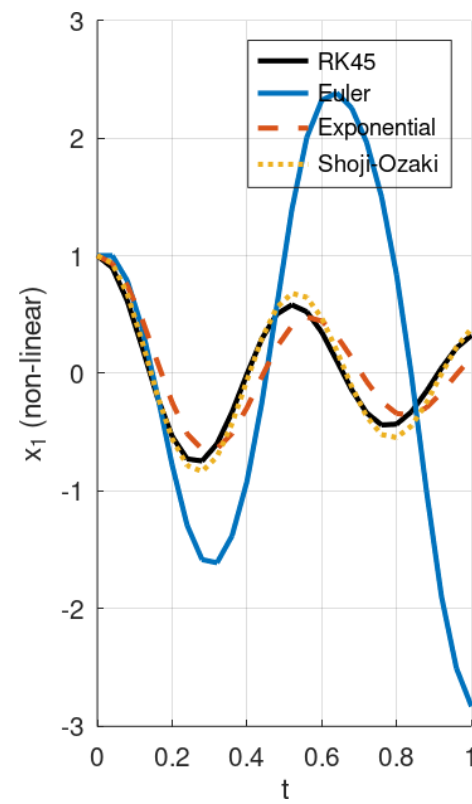
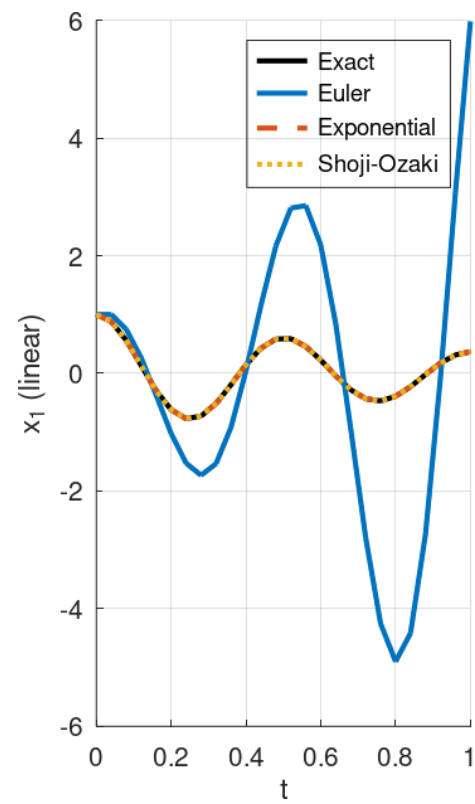
Discretization: example

$$\dot{x} = -3x \quad \dot{x} = -3x - x^2$$



Discretization: example

$$\dot{\mathbf{x}} = \begin{bmatrix} 0 & 1 \\ -\omega^2 & -\gamma \end{bmatrix} \mathbf{x} \quad \dot{\mathbf{x}} = \begin{bmatrix} x_2 \\ -\omega^2 \sin(x_1) - \gamma x_2 \end{bmatrix}$$



Discretization

The system $\dot{\mathbf{x}} - \mathbf{F}\mathbf{x} = \mathbf{L}\mathbf{w}(t)$ has an exact solution between t and $t + \Delta t$

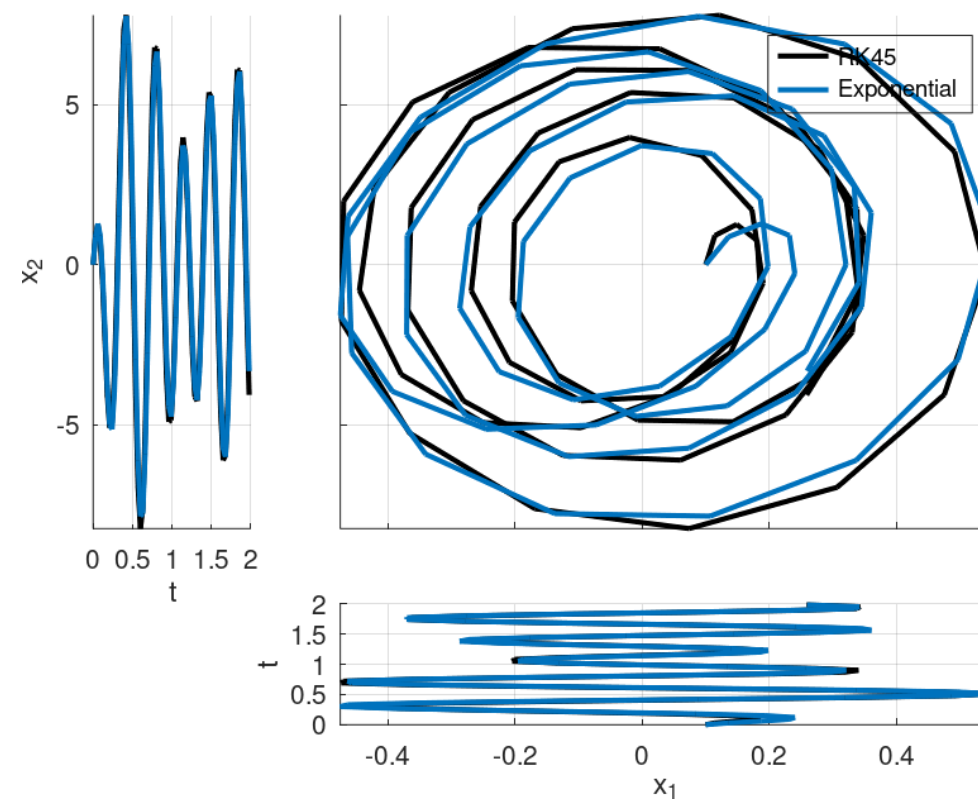
$$\mathbf{x}(t) = \exp(\mathbf{F}\Delta t)\mathbf{x}(0) + \exp(\mathbf{F}\Delta t) \int_t^{t+\Delta t} \exp(\mathbf{F}(t - \tau))\mathbf{L}\mathbf{w}(\tau)d\tau$$

So the approximation is only for the last integral. For example:

$$\hat{\mathbf{x}}(t_{k+1}) = \exp(\mathbf{F}\Delta t) [\hat{\mathbf{x}}(t_k) + \mathbf{L}\mathbf{w}(t_k)\Delta t]$$

When $\mathbf{w}(t)$ is stochastic this is an instance of the Euler-Murayama method.

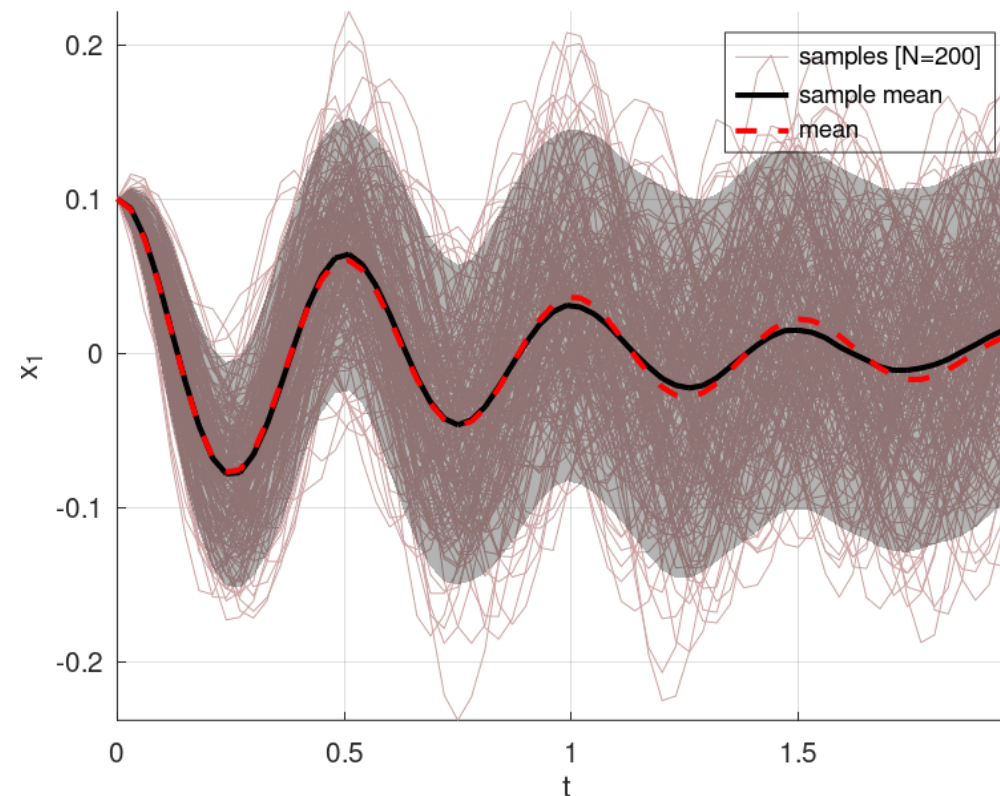
Discretization



source: `s_exp_step_input.m`

Discretization

This also works when $w(t)$ is a random process.



source: `s_exp_step_input.m`

When $w(t)$ is stochastic this is an instance of the Euler-Murayama method.

The stochastic case is subtle, e.g. the discretized noise has a variance that depends on the time step Δt . For details, refer to Simo Särkkä and Arno Solin (2019). Applied Stochastic Differential Equations.

There are stochastic methods that generalize the deterministic ones, e.g. Stochastic Runge-Kutta. For second order system check the Stochastic Verlet algorithm.

Polynomial regression

The data model is $y_k = p(t_k, n) + \epsilon_k$, where

$$p(t, n) := \sum_{s=0}^n a_s t^s$$

is a polynomial of degree n . Consider the system

$$\begin{aligned} \dot{x}_1 &= x_2 & \dot{x}_2 &= x_3 & \dot{x}_3 &= 0 \\ \begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \end{aligned}$$

if we start at $\mathbf{x}(0) = [a_0 \ a_1 \ 2a_2]$ we get

$$x_3(t) = 2a_2 \quad x_2(t) = a_1 + 2a_2 t \quad x_1(t) = a_0 + a_1 t + a_2 t^2$$

we recover $p(t, 2)$ on $x_1(t)$!

Polynomial regression

Can be generalized to any finite degree:

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{F}\mathbf{x}, & F_{ij} &= \delta_{i(j+1)}, \quad 1 \leq i, j \leq n+1 \\ \mathbf{x}(0)^\top &= [a_0 \ a_1 \ 2a_2 \ \dots \ n!a_n]\end{aligned}$$

has $p(t, n)$ in $x_1(t)$.

Then filter

$$\begin{aligned}\mathbf{A} &= \exp(F\Delta t), & \mathbf{Q} &\neq \mathbf{0} \\ \mathbf{H} &= [1 \ 0 \ \dots \ 0], & \mathbf{R} &\neq \mathbf{0}\end{aligned}$$

implements polynomial regression with drift model in the coefficients.

See `s_polyreg.m`

Examples

Multiple Object tracking in 1D

Refer to `tracking_example.py`.

- For more than 2 targets: position information
- For 2 targets: position of 1st, relative position of 2nd w.r.t. 1st

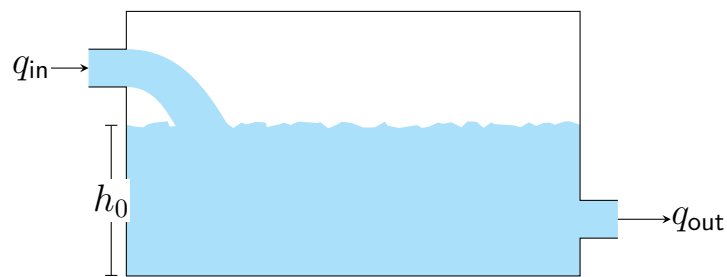
The "true" motion of target n (with mass m_n) is given by:

$$\begin{bmatrix} \dot{p}_n \\ \dot{v}_n \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & \frac{1}{m_n} \end{bmatrix} \begin{bmatrix} p \\ v \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{1}{m_n} \end{bmatrix} f(t)$$

where $f(t)$ is a chaotic force. It is modelled with:

$$\begin{bmatrix} \dot{\tilde{p}}_n \\ \dot{\tilde{v}}_n \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & \frac{1}{m_n} \end{bmatrix} \begin{bmatrix} \tilde{p} \\ \tilde{v} \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \eta(t), \quad \eta(t) \sim \mathcal{N}(0, \sigma_\eta^2)$$

Water tank



The water level in the tank is given by: $\dot{h} = q_{in}(t) - q_{out}(t)$
Discharge is a function of the water level: $q_{out}(h(t)) = Ch^s$, e.g. $s = \frac{1}{2}$.
For any $s \neq 1$, it is a nonlinear system.

For small variations of the level $h(t) = h_o + \Delta h(t)$:

$$\begin{aligned} q_{out}(h(t)) &= Ch(t)^s \simeq C [h_o^s + sh_o^{s-1}(h(t) - h_o)] = C(1 - s)h_o^s + Csh_o^{s-1}h(t) \\ &:= -C_1(h_o) - C_2(h_o)h \\ \dot{h} &= C_2(h_o)h + q_{in}(t) + C_1(h_o) \end{aligned}$$

The linearization adds a constant (negative) term to the input.

Water tank

Linearized dynamics

$$\dot{h} = C_2(h_o)h + q_{\text{in}}(t) + C_1(h_o)$$

The input has a known contribution $q(t)$, and an unknown extra (small) inflow $u(t)$:

$$q_{\text{in}}(t) = q(t) + u(t)$$

We will model it with :

$$\begin{bmatrix} \dot{h} \\ \dot{\delta}_{C_1} \\ \dot{\mathbf{u}} \end{bmatrix} = \begin{bmatrix} C_2(h_o) & 1 & \mathbf{U}(t) \\ 0 & 0 & \mathbf{0} \end{bmatrix} \begin{bmatrix} h \\ \delta_{C_1} \\ \mathbf{u} \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \\ \mathbf{0} \end{bmatrix} [q(t) + C_1(h_o)] + \mathbf{L}\boldsymbol{\eta}(t)$$

$$\dot{\mathbf{x}} = \mathbf{F}(t)\mathbf{x} + \mathbf{B}\tilde{q}(t) + \mathbf{L}\boldsymbol{\eta}(t), \quad \boldsymbol{\eta}(t) \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\eta)$$

where $\mathbf{U}(t)$ is a feature vector (a given set of time-dependent functions), and \mathbf{u} their coefficients. δ_{C_1} is a correction to the constant $C_1(h_o)$ produced by the linearization.

The unknown inflow is modelled as

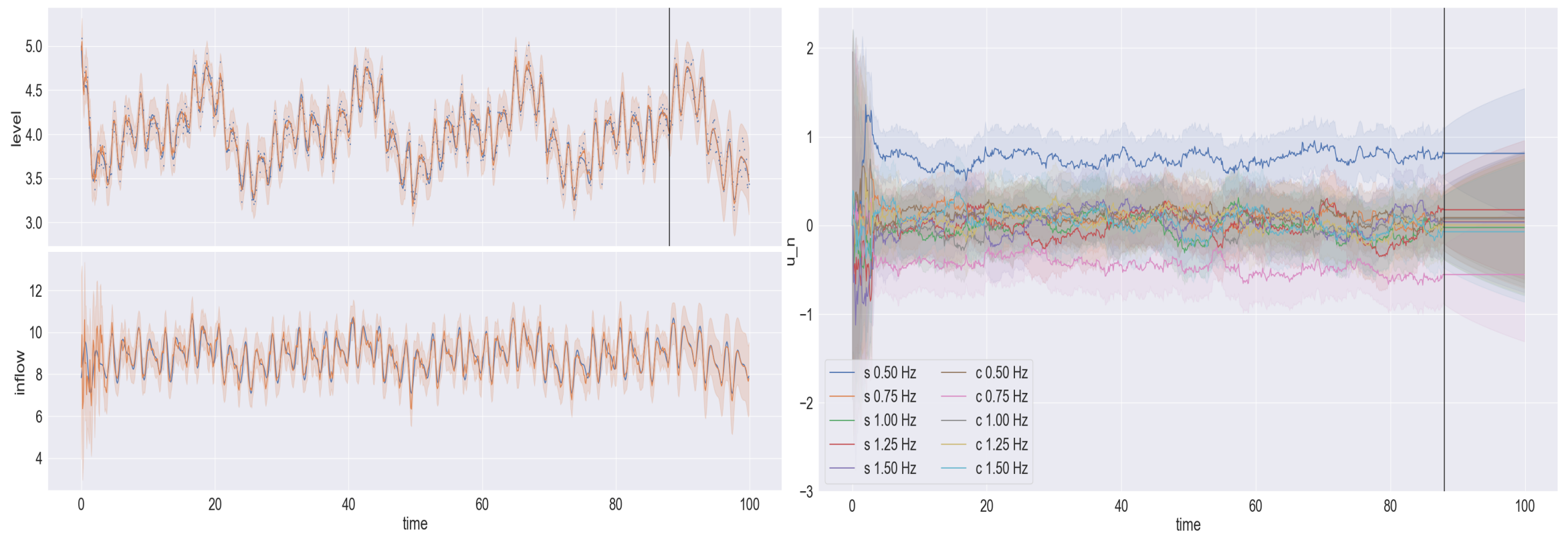
$$\tilde{u}(t) = \sum_{n=1}^N u_n U_n(t)$$

and the known inflow is combined with the constant produced by the linearization

$$\tilde{q}(t) = q(t) + C_1(h_o)$$

Water tank

Refer to `tank_example_basis.py`



Participant examples

Examples provided by participants

- Manuel Regueiro-Picallo (manuel.regueiro@eawag.ch) from Eawag presented the problem of estimating temperature distributions between measurements at two points. He discretized the unsteady diffusion equation $\frac{\partial T}{\partial t} = k \frac{\partial^2 T}{\partial x^2}$ and built the Kalman filter with the resulting matrices.
- Reto Christen (reto.christen@ost.ch) from IET presented the tracking of rapid phase changes in signals from the electrical grid (very high noise!). He also discussed the implementation of the filter in embedded electronics. He used recursive regression with periodic functions of different frequencies.
- Adrian Rohner (adrian.rohner@ost.ch) and Florian Hammer (florian.hammer@ost.ch) from IET discussed the issue of smoothing speed wind velocity measurements from wind turbines, and the evaluation of uncertain speed-to-power curves.

Parameter estimation

What to estimate?

The linear Kalman filter is defined with the parameters \mathbf{A} , \mathbf{H} , \mathbf{Q} , \mathbf{R} , \mathbf{P}_0 , \mathbf{m}_0 .

- \mathbf{A} , \mathbf{H} reflect knowledge of the physical system (unless data-driven).
- Bounds for the diagonal of \mathbf{R} can be estimated from the data.
- \mathbf{Q} in general is hidden, unless derived from physical SDE.
- \mathbf{P}_0 , \mathbf{m}_0 can be critical for the performance of filters with long-term memories (eigenvalues of \mathbf{A} closer to or bigger than 1).

Idea: start by freezing degrees of freedom (à la stochastic gradient decent).

To estimate \mathbf{R} use any smoother on the data and compute the variance of the residuals between smoothed data and raw data. There is almost always a strong expectation on this parameter.

Many physics based models will have limited memory. An exception would be a fluid model where waves can carry information for long distances and induce long-term temporal correlations. In general systems with echo (transport, waves) can produce these effects.

Augmented state

Augment your state with the parameters with constant noisy dynamics (drift model)

$${}^{\text{extra}}x_k = {}^{\text{extra}}x_{k-1} + \epsilon_k, \quad \epsilon_k \sim \mathcal{N}(0, \sigma_{\text{extra}}^2)$$

Works out-of-the-box if resulting model is linear in parameters, i.e. if we have in the dynamics a term of the form ax_{k-1} , adding a to the state will render the model non-linear.

σ_{extra} controls how "reactive" or "nervous" the parameter is. Larger values, quicker adaptation, larger confidence intervals.

See `s_fourier_adaptive.m`

Posterior likelihood

Bayesian filter for state space models with parameter vector θ

$$\begin{aligned}\theta &\sim p(\theta), & \mathbf{x}_0 &\sim p(\mathbf{x}_0|\theta) \\ \mathbf{x}_k &\sim p(\mathbf{x}_k|\mathbf{x}_{k-1}, \theta) \\ \mathbf{y}_k &\sim p(\mathbf{y}_k|\mathbf{x}_k, \theta)\end{aligned}$$

The idea is to compute the posterior distribution of states and parameters given the measurements

$$p(\mathbf{x}_{0:T}, \theta | \mathbf{y}_{1:T}) \propto p(\mathbf{y}_{1:T} | \mathbf{x}_{0:T}, \theta) p(\mathbf{x}_{0:T} | \theta) p(\theta)$$

where

$$\begin{aligned}p(\mathbf{y}_{1:T} | \mathbf{x}_{0:T}, \theta) &= \prod_{k=1}^T p(\mathbf{y}_k | \mathbf{x}_k, \theta) \\ p(\mathbf{x}_{0:T} | \theta) &= p(\mathbf{x}_0 | \theta) \prod_{k=1}^T p(\mathbf{x}_k | \mathbf{x}_{k-1}, \theta)\end{aligned}$$

We focus on the parameters, so we "average" (marginalize) out the states.

These are the same relations as before, we just made the parameters explicit.

$p(\theta)$ is the prior over the parameters.

Posterior likelihood

Bayesian filter for state space models with parameter vector θ

$$\theta \sim p(\theta), \quad x_0 \sim p(x_0|\theta)$$

$$x_k \sim p(x_k|x_{k-1}, \theta)$$

$$y_k \sim p(y_k|x_k, \theta)$$

The posterior distribution of parameters given the measurements

$$p(\theta|y_{1:T}) \propto p(y_{1:T}|\theta)p(\theta)$$

where (check the "averaged" states)

$$p(y_{1:T}|\theta) = \prod_{k=1}^T p(y_k|y_{1:k-1}, \theta)$$

$$p(y_k|y_{1:k-1}, \theta) = \int \underbrace{p(y_k|x_k, \theta)}_{\text{measurement}} \underbrace{p(x_k|y_{1:k-1}, \theta)}_{\text{prediction}} dx_k$$

update

For convenience of the notation we used

$$p(y_1|y_{1:0}, \theta) = p(y_1|\theta)$$

To get the recursion for the posterior of the parameters, the key is to write a recursion for likelihood of the data. We do that next.

Posterior likelihood

We know how to compute the posterior distribution of parameters given the measurements at a given step k

$$p(\boldsymbol{\theta}|\mathbf{y}_{1:k}) \propto p(\mathbf{y}_{1:k}|\boldsymbol{\theta})p(\boldsymbol{\theta})$$

We define log-likelihood $\phi_k(\boldsymbol{\theta})$:

$$\begin{aligned} \exp(\phi_k(\boldsymbol{\theta})) &:= p(\mathbf{y}_{1:k}|\boldsymbol{\theta}) = \prod_{n=1}^k p(\mathbf{y}_n|\mathbf{y}_{1:n-1}, \boldsymbol{\theta}) = p(\mathbf{y}_k|\mathbf{y}_{k-1}, \boldsymbol{\theta}) \underbrace{\prod_{n=1}^{k-1} p(\mathbf{y}_n|\mathbf{y}_{1:n-1}, \boldsymbol{\theta})}_{p(\mathbf{y}_{1:k-1}|\boldsymbol{\theta}) = \exp(\phi_{k-1}(\boldsymbol{\theta}))} \\ &= p(\mathbf{y}_k|\mathbf{y}_{k-1}, \boldsymbol{\theta}) \exp(\phi_{k-1}(\boldsymbol{\theta})) \end{aligned}$$

Taking logs

$$\phi_k(\boldsymbol{\theta}) = \phi_{k-1}(\boldsymbol{\theta}) + \log [p(\mathbf{y}_k|\mathbf{y}_{k-1}, \boldsymbol{\theta})]$$

with $\phi_0(\boldsymbol{\theta}) = \log [p(\boldsymbol{\theta})]$

If the prior $p(\boldsymbol{\theta})$ is uniform, this is the same as maximum likelihood (ML). If the distribution reflect our prior knowledge about the parameters, then it is maximum a posteriori probability (MAP).

$\phi_k(\boldsymbol{\theta})$ is the log-likelihood, sometimes the negative is defined.

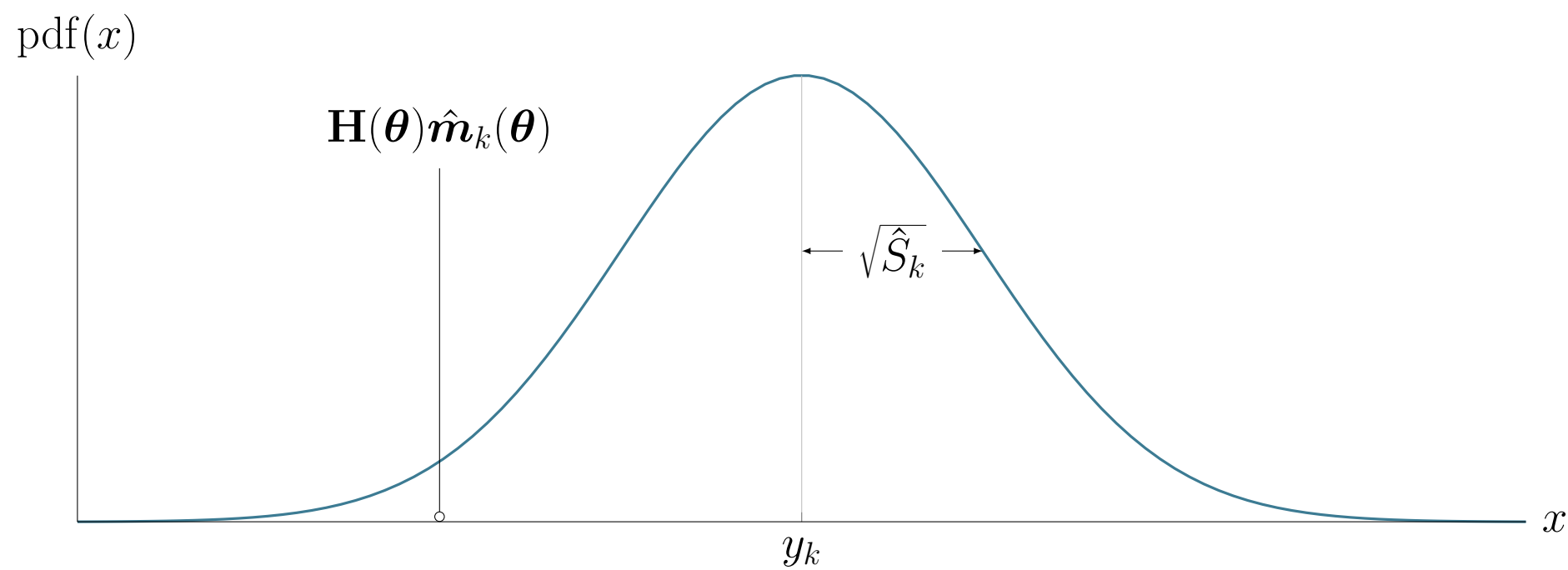
Maximizing the log-likelihood gives optimal estimates. This should remind you of Viterbi, if you know it.

Linear KF likelihood

With the filters we have been studying, the log-likelihood has an explicit formula.

$$\mathbf{v}_k(\boldsymbol{\theta}) := \mathbf{y}_k - \mathbf{H}(\boldsymbol{\theta})\hat{\mathbf{m}}_k(\boldsymbol{\theta})$$

$$\phi_k(\boldsymbol{\theta}) = \phi_{k-1}(\boldsymbol{\theta}) - \frac{1}{2} \log [2\pi \det \hat{\mathbf{S}}_k(\boldsymbol{\theta})] - \frac{1}{2} \mathbf{v}_k^\top(\boldsymbol{\theta}) \hat{\mathbf{S}}_k^{-1}(\boldsymbol{\theta}) \mathbf{v}_k(\boldsymbol{\theta})$$



Maximum likelihood

$$\begin{aligned}\mathbf{x}_k &= \exp\left(\begin{bmatrix} 0 & 1 \\ -\omega^2 & 0 \end{bmatrix} \Delta t\right) \mathbf{x}_{k-1} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} w_{k-1}, & w_{k-1} &\sim \mathcal{N}(0, q) \\ y_k &= [1 \ 0] \mathbf{x}_k + r_k, & r_k &\sim \mathcal{N}(0, \sigma_y^2)\end{aligned}$$

`s_estimate_oscillator.py` does

- Line search for ω
- Optimization for $\omega, q, \sigma_y, \mathbf{m}_0$

Gradient based methods will not be sensitive to the data if started far away from a local maximum: likelihood tends to be sharp and wiggly, vanishing gradients. Alternative: informative priors, e.g. tight bounds on the parameters.

If you want to optimize many parameters of the model, check the Expectation Maximization algorithm.

Wrap-up

Summary and Outlook

- Linear Kalman filter with Gaussian distributions ✓
- Non-linear
 - Local linearization (Extended Kalman filter)
 - Gaussian approximation (Unscented Kalman filter, General Gaussian filters: GHKF, CKF)
 - Particle filters, ...
 - Brute force (naive MC, simulation of trajectories)
- Other noise models (check generalizations of KF, Särkkä's book references)
- Recursive or Batch?
 - Data: fixed size or incoming?
 - Estimation parameters: recorded data (batch) or online?