

Filtering and Smoothing

In quite many applications, we are concerned with estimating the latent variables from their observations / measurements. For example, estimating the trajectory of moving object based on its sensors readings (e.g., gyro and acceleration, angular velocity); estimating COVID cases from city waste water.

In our context, we model the latent variable by an SDE. More specifically, we consider a model

$$dX(t) = a(X(t)) dt + b(X(t)) dW(t), \quad X(0) = X_0,$$

$$Y_k | X(t_k) \sim P_{Y_k | X(t_k)}(Y_k | X_k) \quad \text{We will abbreviate this as } P(Y_k | X_k)$$

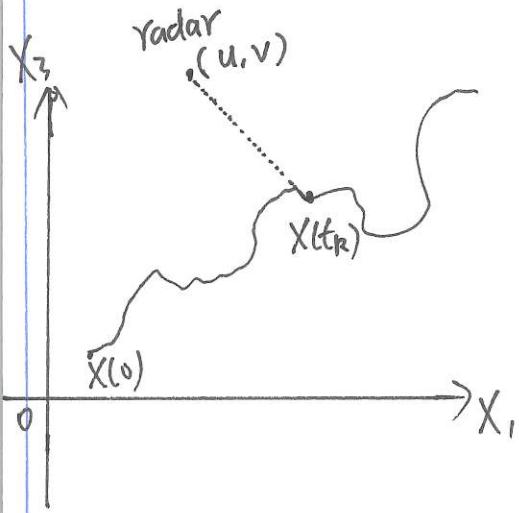
where the process \underline{X} stands for the latent/unobserved process, and the random variable Y_k represents the ~~actual~~ measurements at time t_k , which follows a conditional probability density function $P(Y_k | X_k)$ given.

Example 28. Motion tracking model - 2D.

$$d \begin{bmatrix} X_1(t) \\ X_2(t) \\ X_3(t) \\ X_4(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} X_1(t) \\ X_2(t) \\ X_3(t) \\ X_4(t) \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix} dW(t)$$

$w \in \mathbb{R}^2$

$$Y_k = \begin{bmatrix} \sqrt{(X_1(t_k) - u)^2 + (X_3(t_k) - v)^2} \\ \arctan\left(\frac{X_3(t_k) - v}{X_1(t_k) - u}\right) \end{bmatrix} + \zeta_k, \quad \zeta_k \sim \text{N}(0, I_2)$$



In this example, X_1 and X_3 represent the abscissa and ordinate of the object, respectively. The corresponding velocities are X_2 and X_4 . We assume that there is a sensor/radar placed at (u, v) that measures the relative distance and bearing of the object. We would like to track the object's state based on the measurements.

→ Tracking in what sense?

T3

There are ~~as~~ various notions in estimating the latent process from the measurements. Recall the tracking example in Example 28, ~~we would like to~~ and suppose that at time t_{k-1} we have an estimate of the object. We would like to estimate the state of the object at time t_k based on the previous estimate at t_{k-1} and the measurement at the present time t_k . This gives the heuristics of the filtering problem, and this is formulated as follows.

Filtering

Suppose that we have a sequence of measurements $Y_{1:T} := \{Y_1, Y_2, \dots, Y_T\}$ at times $t_1 < t_2 < \dots < t_T$. The goal of the filtering is to estimate $X(t_k)$ progressively and sequentially for $k=1, 2, \dots$ from the measurements. More precisely, we would like to obtain the filtering probability density function

$$P(X_k | y_{1:k}), \quad \text{shorthand of } P(x(t_k) | Y_{1:k} (x_k | y_{1:k}))$$

which represents the conditional distribution of $X(t_k)$ conditionally on the measurements up to time t_k . ~~Moreover, Furthermore,~~ suppose that we know the previous filtering density $P(X_{k-1}|Y_{1:k-1})$ we want to express $P(X_k|Y_{1:k})$ in terms of $P(X_{k-1}|Y_{1:k-1})$, so that the filtering algorithm is an online recursive routine

$$P(X_0) \rightarrow P(X_1|Y_{1:1}) \rightarrow P(X_2|Y_{1:2}) \rightarrow \dots P(X_k|Y_{1:k}) \dots$$

We now show how to do so by the Markov property and Baye's rule.

By Baye's rule

$$\begin{aligned} P(X_k|Y_{1:k}) &= P(X_k|Y_{1:k-1}, Y_k) = \frac{P(Y_k|X_k) P(X_k|Y_{1:k-1})}{P(Y_k|Y_{1:k-1})} \\ &= \frac{P(Y_k|X_k) \underbrace{P(X_k|Y_{1:k-1})}_{\int P(Y_k|X_k) P(X_k|Y_{1:k-1}) dX_k}}{\int P(Y_k|X_k) P(X_k|Y_{1:k-1}) dX_k}, \end{aligned}$$

where $\xrightarrow{\text{predictive density}}$ $P(X_k|Y_{1:k-1})$ $\xrightarrow{\text{transition density}}$ $\frac{P(X_k|X_{k-1})}{P(X_{k-1}|Y_{1:k-1})} P(X_{k-1}|Y_{1:k-1})$ $\xrightarrow{\text{the previous filtering density}}$

It is then clear that we can use $P(X_{k-1}|Y_{1:k-1})$ to compute $P(X_k|Y_{1:k-1})$ then use ~~the prediction~~ correct the prediction $P(X_k|Y_{1:k-1})$ based on $P(Y_k|X_k)$ to yield $P(X_k|Y_{1:k})$.

Algorithm 29 Filter

Input : $y_{1:T}$, and initial condition $p(x_0)$

Define

$$p(x_0 | y_{1:0}) := p(x_0)$$

Output : $p(x_1 | y_{1:1}), p(x_2 | y_{1:2}), \dots p(x_T | y_{1:T})$

1. For $k=1, 2, \dots T$ do

$$p(x_k | y_{1:k-1}) = \int p(x_k | x_{k-1}) p(x_{k-1} | y_{1:k-1}) dx_{k-1}$$

$$3. p(x_k | y_{1:k}) = \frac{p(y_k | x_k) p(x_k | y_{1:k-1})}{\int p(y_k | x_k) p(x_k | y_{1:k-1}) dx_k}$$

Return $p(x_1 | y_{1:1}), p(x_2 | y_{1:2}), \dots p(x_T | y_{1:T})$ $O(T)$ complexity

There is also another notion of the estimation, called smoothing, which is widely used as well especially in machine learning. The smoothing refers to estimate the state based on all the measurements unlike the filtering which uses only the historical measurements.

More specifically, the smoothing density is

$$p(x_k | y_{1:T}), \text{ for } k=1, 2, \dots T.$$

and we can compute it by

Previous smoothing density at t_{k+1}

$$p(x_k | y_{1:T}) = \int p(x_k | x_{k+1}, y_{1:T}) \underbrace{p(x_{k+1} | y_{1:T})}_{\text{Previous smoothing density at } t_{k+1}} dx_{k+1}$$

Thanks to the Markov property, we have

$$\begin{aligned} p(x_k | x_{k+1}, y_{1:T}) &= p(x_k | x_{k+1}, y_{1:k}) \\ &= \frac{p(x_{k+1} | x_k, y_{1:k}) p(x_k | y_{1:k})}{p(x_{k+1} | y_{1:k})} \\ &= \frac{p(x_{k+1} | x_k) p(x_k | y_{1:k})}{p(x_{k+1} | y_{1:k})} \end{aligned}$$

Substitute the expression for $p(x_k | x_{k+1}, y_{1:T})$ back to $p(x_k | y_{1:T})$ we arrive at

$$p(x_k | y_{1:T}) = p(x_k | y_{1:k}) \int \frac{p(x_{k+1} | x_k) p(x_{k+1} | y_{1:T})}{p(x_{k+1} | y_{1:k})} dx_{k+1}.$$

Hence, to obtain the smoothing density at any $t_k \leq t_i$, we can first run a filter forward to get $p(x_1 | y_1)$, $p(x_2 | y_{1:2})$, ... $p(x_T | y_{1:T})$, then backward compute $p(x_{T-1} | y_{1:T})$, $p(x_{T-2} | y_{1:T})$, ... $p(x_{k+1} | y_{1:T})$, $p(x_k | y_{1:T})$. This is summarised in the following algorithm.

XX

Algorithm 29. Smoothing

Input: $y_{1:T}$, $p(x_0)$

Output: $p(x_1 | y_{1:T})$, $p(x_2 | y_{1:T})$, ... $p(x_T | y_{1:T})$

1. Run the filter to obtain

$$p(x_1 | y_{1:T}), p(x_2 | y_{1:T}), \dots, p(x_T | y_{1:T})$$

2. For $k = T-1, T-2, \dots, 1$ do

$$p(x_{k|T} | y_{1:T}) = p(x_{k|1} | y_{1:k}) \frac{\int p(x_{k+1} | x_k) p(x_{k+1} | y_{1:T}) dx_{k+1}}{p(x_{k+1} | y_{1:k})}$$

Return $p(x_1 | y_{1:T}), p(x_2 | y_{1:T}), \dots, p(x_T | y_{1:T})$

Unfortunately, the filtering and smoothing algorithms are in reality not implementable,

Remark: the predictive density is already computed in the filtering routine. No need to compute it again.

because 1) the transition density $p(x_k | x_{k-1})$ is intractable.
 2) the computations for the integrals are intractable (e.g., $\int p(y_k | x_k) p(x_k | y_{1:k-1}) dx_k$). Therefore, in practice we have to approximate the filtering and smoothing densities, except for a few isolated cases, such as linear Gaussian systems. ~~In what~~

In what follows, we introduce the celebrated Kalman filter and Rauch-Tung-Striebel smoother.

RTS

↳ Kalman filter and RTS smoother

The Kalman filters and RTS smoothers ~~works~~ work on linear SDEs and linear Gaussian measurement models.

More specifically, the models are of the form

$$dX(t) = A X(t) dt + B dW(t),$$

$$Y_k = H X(t_k) + \varepsilon_k, \quad \varepsilon_k \sim N(0, \Sigma_k),$$

where $X \in \mathbb{R}^d$, $A \in \mathbb{R}^{d \times d}$, $B \in \mathbb{R}^{d \times n}$, $H \in \mathbb{R}^{s \times d}$, $\varepsilon_k \in \mathbb{R}^s$. Of course we can also let the SDE coefficients A and B be time-dependent (see, Lecture 5), but for simplicity, let us consider them time-invariant.

As we have seen in Lecture 5 that solution to the linear SDE is Gaussian process, hence, ~~*~~ the measurement random variables Y_k 's are Gaussian too. Consequently, all the

PDFs in the filtering and smoothing algorithms are Gaussian. Hence, for this model, we can solve the filtering and smoothing problems in closed-form by playing with "Gaussian algebras". The most important algebra we are going to use is the Gaussian identity.

Remark 30. Gaussian identity. Let $\alpha \in \mathbb{R}^{d\alpha}$ and $\beta \in \mathbb{R}^{d\beta}$ be two jointly Normal distributed vectors, such that

$$\begin{bmatrix} \alpha \\ \beta \end{bmatrix} \sim N\left(\begin{bmatrix} \mathbb{E}[\alpha] \\ \mathbb{E}[\beta] \end{bmatrix}, \begin{bmatrix} \text{Cov}[\alpha], \text{Cov}[\alpha, \beta] \\ \text{Cov}[\beta, \alpha] \text{Cov}[\beta] \end{bmatrix}\right)$$

Then the distribution of α condition on β is also Normal with mean

$$\mathbb{E}[\alpha | \beta] = \mathbb{E}[\alpha] + \text{Cov}[\alpha, \beta] (\text{Cov}[\beta])^{-1} (\beta - \mathbb{E}[\beta])$$

and covariance

$$\text{Cov}[\alpha | \beta] = \text{Cov}[\alpha] - \text{Cov}[\alpha, \beta] (\text{Cov}[\beta])^{-1} \text{Cov}[\beta, \alpha].$$

~~Mutatis mutantis~~ for β condition on α .

~~Mutatis mutandis~~

We are now ready to derive the filtering densities. Since the they are Gaussian, denote $P(X_k | Y_{1:k}) := N(\bar{X}_k | M_k, V_k)$, where M_k and V_k stand for the filtering mean and covariance, respectively. ~~Also~~ Also remark that M_k and V_k depend on the data $Y_{1:k}$. ~~since~~ The first step in the filtering algorithm 29 requires to compute

$$P(X_k | Y_{1:k-1}) = \int P(X_k | X_{k-1}) P(X_{k-1} | Y_{1:k-1}) dX_{k-1}.$$

The properties of the linear SDE tell us that

$$\begin{aligned} P(X_k | X_{k-1}) &= N(X_k | \mathbb{E}[X(t_k) | X(t_{k-1})], \text{Cov}[X(t_k) | X(t_{k-1})]) \\ &= N(X_k | F_k X_{k-1}, \Sigma_k). \end{aligned}$$

Please recall how to compute F_k and Σ_k in Lecture 5.

Hence,

$$\begin{aligned} P(X_k | Y_{1:k-1}) &= \int N(X_k | F_k X_{k-1}, \Sigma_k) N(X_{k-1} | M_{k-1}, V_{k-1}) dX_{k-1} \\ &= N(X_k | F_k M_{k-1}, F_k V_{k-1} F_k^T + \Sigma_k) \\ &:= N(X_k | \bar{M}_k, \bar{V}_k) \quad \text{shorthand notation} \end{aligned}$$

Next, we apply the Gaussian identity in Remark 30 to obtain $p(x_k | y_{1:k})$. Now imagine that the α and β in Remark 30 are $x_k | y_{1:k-1}$ and $y_k | y_{1:k-1}$, respectively, then

$$m_k := \mathbb{E}[x_k | y_{1:k}] = \mathbb{E}[x_k | y_{1:k-1}] + \text{Cov}[x_k, y_k | y_{1:k-1}] \cdot \\ \text{Cov}[y_k | y_{1:k-1}]^{-1} (y_k - \mathbb{E}[y_k | y_{1:k-1}]),$$

$$v_k := \text{Cov}[x_k | y_{1:k}] = \text{Cov}[x_k | y_{1:k-1}] - \text{Cov}[x_k, y_k | y_{1:k-1}] \cdot \\ \text{Cov}[y_k | y_{1:k-1}]^{-1} \text{Cov}[y_k, x_k | y_{1:k-1}].$$

Recall that

$$p(y_k | y_{1:k-1}) = \int p(y_k | x_k) p(x_k | y_{1:k-1}) dx_k \\ = \int N(y_k | Hx_k, \Sigma_k) N(x_k | m_k^-, V_k^-) dx_k \\ = N(y_k | Hm_k^-, HV_k^-H^T + \Sigma_k)$$

and that

$$\text{Cov}[x_k, y_k | y_{1:k-1}] = \text{Cov}[x_k, Hx_k + \eta_k | y_{1:k-1}] \\ = V_k^- H^T.$$

Putting everything together we have

$$m_k = \bar{m}_k + V_k^{-} H^T (H V_k^{-} H^T + \Sigma_k)^{-1} (y_k - H \bar{m}_k)$$

$$V_k = V_k^{-} - V_k^{-} H^T (H V_k^{-} H^T + \Sigma_k)^{-1} H V_k^{-}.$$

We ~~can~~ summarise the Kalman filtering in the following algorithm.

Algorithm³ Kalman filter

Input : Measurements $y_{1:T}$, initial mean m_0 and covariance V_0

Output: $\{m_k, V_k\}_{k=1}^T$

1. For $k = 1, 2, \dots, T$ do

2. Compute F_k and Σ_k // can be pre-computed

3. $\bar{m}_k = F_k m_{k-1}$ > prediction

4. $\bar{V}_k = F_k V_{k-1} F_k^T + \Sigma_k$

5. $S_k = H_k V_k^{-} H_k^T + \Sigma_k$

6. Kalman gain $K_k = V_k^{-} H_k^T S_k^{-1}$ update

7. $m_k = \bar{m}_k + K_k (y_k - H_k \bar{m}_k)$

8. $V_k = V_k^{-} - K_k S_k K_k^T$

Return $\{m_k, V_k\}_{k=1}^T$

Note that V_k does Not depend on the data values $y_{1:k}$.

By using a similar derivation, we can obtain the RTS smoother for the smoothing densities. This is left as a homework for the readers. The RTS smoother is summarised as follows.

$$\text{Denote } P(x_k | y_{1:T}) \\ := N(x_k | m_k^s, V_k^s)$$

Algorithm 82. RTS smoother

Input: Measurements $y_{1:T}$, initial mean m_0 and covariance V_0 .

Output: $\{m_k^s, V_k^s\}_{k=1}^T$

1. Do the Kalman filter to obtain

$$\{m_k, V_k\}_{k=1}^T$$

2. Let $m_T^s = m_T$ and $V_T^s = V_T$

3. For $k = T-1, T-2, \dots, 1$ do

$$\bar{m}_{k+1} = F_{k+1} m_k$$

$$\bar{V}_{k+1} = F_{k+1} V_k F_{k+1}^T + \Sigma_{k+1}$$

$$G_k = V_k F_{k+1}^T (\bar{V}_{k+1})^{-1} \quad \text{smoother gain}$$

$$m_k^s = m_k + G_k (m_{k+1}^s - \bar{m}_{k+1})$$

$$V_k^s = V_k + G_k (V_{k+1}^s - \bar{V}_{k+1}) G_k^T$$

Return $\{m_k^s, V_k^s\}_{k=1}^T$

these two predictive quantities are already computed in the filter.
no need to compute them again.

Example.33. Consider a simple motion model

$$d \begin{bmatrix} X_1(t) \\ X_2(t) \\ X_3(t) \\ X_4(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} X_1(t) \\ X_2(t) \\ X_3(t) \\ X_4(t) \end{bmatrix} dt + \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix} dw(t)$$

$$Y_k = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} X(t_k) + h_k, \quad h_k \sim N(0, \Sigma)$$

$$X(0) \sim N(0, I_4).$$

1) Generate a pair of trajectory and measurements $(X_{1:T}, Y_{1:T})$ at times t_1, t_2, \dots, t_T

2) use Kalman filter and RTS smoother to estimate $X_{1:T}$ from $Y_{1:T}$. That is, compute the filtering and smoothing densities at the times.

See, 'lec6_kalman_filter_rts_smoother.ipynb' for how to do so.

↳ Parameter estimation with maximum likelihood

In Lecture 4, we have already talked about the parameter estimation of SDEs. However, we assumed that we can directly observe the SDE. But, in reality, the SDE is usually a latent process, ~~and~~ which we do not directly observe. Consider a model as follows.

$$dX(t) = a(X(t); \theta) dt + b(X(t); \theta) dW(t), \quad X(0) = X_0,$$
$$Y_k | X(t_k) \sim P(Y_k | X_{k-1}; \theta).$$

where $\theta \in \mathbb{R}^d$ is an unknown parameter that we want to estimate from the measurements $\{y_1, y_2, \dots, y_T\} := \mathcal{Y}_{1:T}$

$\mathcal{Y}_{1:T} := \{Y_1, Y_2, \dots, Y_T\}$. To estimate θ , we can use the maximum likelihood method, and the likelihood is

$$P(\mathcal{Y}_{1:T}; \theta) = \int \prod_{k=1}^T P(Y_k | X_k; \theta) \prod_{k=1}^T P(X_k | X_{k-1}; \theta) P(X_0) dX_{1:T}. \quad (1)$$

Remark that compared to the parameter estimation in Lecture 4, we now need to integrate out the variables $X_{1:T}$ in order to

Compute $P(y_{1:T}; \theta)$ which is way harder than that, and is not analytically tractable in general, because of the intractable transition density $P(x_k | x_{k-1})$.

The equation in ① is not the only way to express the likelihood. We can also factorise $P(y_{1:T}; \theta)$ in terms of the filtering densities :

$$\begin{aligned} P(y_{1:T}; \theta) &= \prod_{k=1}^T P(y_k | y_{1:k-1}; \theta) \\ &= \prod_{k=1}^T \underbrace{\int P(y_k | x_k; \theta) P(x_k | y_{1:k-1}; \theta) dx_k}_{\downarrow} \\ &= \int P(x_k | x_{k-1}; \theta) P(x_{k-1} | y_{1:k-1}) dx_{k-1}. \end{aligned}$$

Hence we can compute $P(y_{1:T}; \theta)$ along with the run of the filter.

Example 34.

$$dx(t) = A(\theta) x(t) dt + B(\theta) dW(t)$$

$$x_k = H(\theta) x(t_k) + \varepsilon_k, \quad \varepsilon_k \sim N(0, \Sigma).$$

$$\begin{aligned}
 P(Y_{1:T}; \theta) &= \prod_{k=1}^T \int p(Y_k | X_k; \theta) p(X_k | \cancel{Y_{1:k-1}}; \theta) dX_k \\
 &= \prod_{k=1}^T \int N(Y_k | H(\theta)X_k, \Sigma) N(X_k | \bar{m}_k(\theta), V_k(\theta)) dX_k \\
 &= \prod_{k=1}^T N(Y_k | H(\theta)\bar{m}_k(\theta), H(\theta)V_k(\theta)^T + \Sigma).
 \end{aligned}$$

Hence the objective function for carrying out the MLE is

$$-\log P(Y_{1:T}; \theta) = \frac{1}{2} \sum_{k=1}^T \left((y_k - H(\theta)\bar{m}_k(\theta))^T S_k(\theta)^{-1} (y_k - H(\theta)\bar{m}_k(\theta)) + \log \det(2\pi S_k(\theta)) \right).$$

Algorithm 35. Computing the negative log likelihood with Kalman filter

Input: The same as with Algorithm 31.

Output: The same as with Algorithm 31.

1. The same as with Algorithm 31, except that
additionally
 we compute

$$nll_k = \frac{1}{2} ((y_k - H_k \bar{m}_k)^T S_k^{-1} (y_k - H_k \bar{m}_k) + \log \det(2\pi S_k))$$

at the end of each loop.

Return $\sum_{k=1}^T nll_k$.

88

Non-linear filtering and smoothing

In the previous section, we have seen how to use Kalman filters and smoothers which are optimal in the least square sense for linear Gaussian models.

However, in reality, we often have to tackle non-linear models. Consider models of the form:

$$dx(t) = a(x(t))dt + b(x(t))dW(t), \quad x(0) = X_0 \sim P(X_0)$$

$$Y_k \mid X(t_k) \sim P(Y_k \mid X_k).$$

It is known challenging to implement Algorithms 28 and 29 to solve the filtering and smoothing densities $P(X_k \mid Y_{1:k})$ and $P(X_k \mid Y_{1:T})$ in closed-form. Hence, we may have to approximate them. The most celebrated approximations are arguably the

Note: if the measurement model is continuous in time, e.g.

$$dy(t) = \int_0^t h(x(s))ds + \int_0^t g(s)dW(s),$$
 the filtering density is governed by

the Kushner-Stratonovich equation which is a stochastic partial differential equation. We can simulate such an SPDE to solve the filtering problem in continuous-time.

Gaussian filters and smoothers, and the particle filters and smoothers. The Gaussian filters include, for instance, the extended \mathcal{K} . For simplicity, in this lecture, we introduce the Gaussian filters and smoothers.

Remark: the backbone of particle filters and smoothers is the sequential Monte Carlo techniques. Essentially, instead of trying to approximate the filtering density $P(X_k | Y_{1:k})$, the particle filters approximate any expectation $E[\varphi(X_{t:n}) | Y_{1:k}]$ for arbitrary function φ by Monte Carlo. You can also think that $P(X_k | Y_{1:k})$ is represented by a bunch of weighted samples. See, Chopin and Papaspiliopoulos, 2020.

§ Gaussian approximate filtering and smoothing
please avoid calling them Gaussian assumed-density filters and smoothers, since you cannot assume something that is not true.
The idea essence of the approximate Gaussian filters and smoothers consists in the following approximations:

- 1) $P(X_k | Y_{1:k}) \approx N(X_k | m_k, V_k)$ Gaussian filtering density
- 2) $P(X_k | Y_{1:T}) \approx N(X_k | m_k^S, V_k^S)$ Gaussian smoothing density

$$3) P(X_k | X_{k-1}) \sim N(X_k | \hat{E}[X(t_k) | X(t_{k-1})], \text{Cov}[X(t_k) | X(t_{k-1})])$$

Gaussian transition

4) $p(x_0) \sim N(x_0 | m_0, V_0)$ Gaussian initial

$$5) P(y_k | X_k) \sim N(y_k | h(x_k), \Sigma_k) \quad y_k = h(x_{(k)}) + h_k, \quad h_k \sim N(0, \Sigma_h)$$

By applying these Gaussian approximations, we can solve the filtering and smoothing problems in closed-form by playing with Gaussian algebras. Substituting these approximations back into the general filtering and smoothing equations in Algorithms 28 and 29, and use a similar derivation for the Kalman filter and RTS smoother, we have the ~~skeleton~~ of Gaussian approximate filter and smoother.

Algorithm³⁶ Gaussian approximate filter.

Input: Measurements $y_{1:T}$, initial m_0 and v_0

output: $\{M_k, V_k\}_{k=1}^T$

1. | For $k = 1, 2, \dots, T$ do .

$$\textcircled{1} \quad \left\{ \bar{m}_k \equiv \int x_k \cdot p(x_k | y_{1:k-1}) dx_k \right.$$

$$\hat{V}_{k|k} \doteq \int (X_{k|k} - M_{k|k}) (X_{k|k} - M_{k|k})^T p(X_{k|k} | y_{1:k+1}) dX_{k|k}$$

Predictions

9 /

update

(2)

$$\left\{ \begin{array}{l} M_k = \int h(\bar{x}_k) N(\bar{x}_k | \bar{m}_k, \bar{v}_k) d\bar{x}_k \\ S_k = \int (h(\bar{x}_k) - M_k)(h(\bar{x}_k) - M_k)^T N(\bar{x}_k | \bar{m}_k, \bar{v}_k) d\bar{x}_k \\ K_k = \left(\int (\bar{x}_k - \bar{m}_k)(\bar{x}_k - \bar{m}_k)^T d\bar{x}_k \right) S_k^{-1} \\ \bar{m}_k = \bar{m}_k + K_k(y_k - M_k) \\ \bar{v}_k = \bar{v}_k + K_k S_k K_k^T \end{array} \right.$$

Return $\{\bar{m}_k, \bar{v}_k\}_{k=1}^T$

Algorithm 37 Gaussian approximate smoothing.

Input: Measurements $y_{1:T}$, initial M_0 and V_0

Output: $\{\bar{m}_k^S, \bar{v}_k^S\}_{k=1}^T$

1. Do the Gaussian filtering to obtain

$\{\bar{m}_k, \bar{v}_k\}_{k=1}^T$ and $\{\bar{m}_k^-, \bar{v}_k^-\}_{k=1}^T$

2. Let $M_T^S = M_T$ and $V_T^S = V_T$

3. For $k = T-1, T-2, \dots, 1$ do

$D_{k+1} \leftarrow \text{Cov}[X_k, X_{k+1} | y_{1:k}]$.

$$D_{k+1} = \int \bar{x}_k \mathbb{E}[X_{k+1} | \bar{x}_k] N(\bar{x}_k | \bar{m}_k, \bar{v}_k) d\bar{x}_k - \bar{m}_k(\bar{m}_{k+1}^-)^T$$

$$G_k = D_{k+1} (\bar{v}_{k+1}^-)^{-1}$$

$$M_k^S = M_k + G_k (M_{k+1}^S - \bar{m}_{k+1}^-)$$

$$V_k^S = V_k + G_k (V_{k+1}^S - \bar{v}_{k+1}^-) G_k^T$$

Return $\{M_k^S, V_k^S\}_{k=1}^T$

Q2

The Gaussian filter and smoother look very similar to the Kalman filter and smoother, except that now we have to solve some Gaussian expectations with non-linear integrands.

These expectations specifically refer to :

$$\textcircled{1} \quad \left\{ \begin{array}{l} \bar{m}_k \approx \int x_k P(x_k | y_{1:k-1}) dx_k \\ V_k \approx \int (x_k - \bar{m}_k)(x_k - \bar{m}_k)^T P(x_k | y_{1:k-1}) dx_k \end{array} \right.$$

$$\textcircled{2} : m_k, S_k, K_k$$

$$\textcircled{3} : D_{k+1} = \int x_{k+1} \mathbb{E}[x_{k+1} | x_k] N(x_k | m_k, V_k) dx_k - m_k (m_{k+1})^T$$

How do we compute the predictive mean and covariance in $\textcircled{1}$? ~~Recall~~ There are a plenty of ways to compute them.

For example, recall Lecture 4, the statistics of SDES, we have

$$\left\{ \begin{array}{l} \frac{d \mathbb{E}[x(t) | y_{1:k-1}]}{dt} = \mathbb{E}[a(x(t)) | y_{1:k-1}] \\ \text{denote } \bar{m}(t) := \mathbb{E}[x(t) | y_{1:k-1}] \quad \text{shorthand } t \in [t_{k-1}, t_k] \\ \frac{d \text{Cov}[x(t) | y_{1:k-1}]}{dt} = \mathbb{E}[a(x(t))(x(t) - \bar{m}(t))^T + (x(t) - \bar{m}(t))a(x(t))^T \\ \quad + b(x(t))b(x(t))^T | y_{1:k-1}] \end{array} \right.$$

Hence, to obtain \bar{m}_k and V_k , we solve the ODE above at t_k starting from the initial values m_k and V_k at t_{k-1} . But how

do we solve the expectation on the RHS of the ODEs?
We can approximate the expectation by Gaussian, then

$$\mathbb{E}[a(x(t)) \mid y_{1:k-1}] \approx \int a(x(t)) N(x(t) \mid \hat{x}^m(t), \hat{V}(t)) dx(t)$$

We can further approximate a by Linearisation

$$a(x) \approx a(\hat{m}) + \boxed{Ja(\hat{m})}(x - \hat{m}) + \dots$$

Jacobian of a at \hat{m}

then we have

$$\begin{aligned} \mathbb{E}[a(x(t)) \mid y_{1:k-1}] &\approx \mathbb{E}[a(\hat{m}(t)) + Ja(\hat{m}(t))(x(t) - \hat{m}(t))] \\ &= a(\hat{m}(t)) \end{aligned}$$

~~$$\mathbb{E}[a(x(t))(x(t) - \hat{m}(t))^T + (x(t) - \hat{m}(t))a(x(t))^T]$$~~

$$+ b(x(t))b(x(t))^T \mid y_{1:k-1}]$$

$$\approx \hat{V}(t) Ja(\hat{m}(t))^T + Ja(\hat{m}(t)) \hat{V}(t) + b(\hat{m}(t))b(\hat{m}(t))^T.$$

putting everything together, we have

$$\left\{ \frac{d\hat{m}(t)}{dt} = a(\hat{m}(t)) \quad , \quad t \in [t_{k-1}, t_k]. \right.$$

$$\left. \frac{d\hat{V}(t)}{dt} = \hat{V}(t) Ja(\hat{m}(t))^T + Ja(\hat{m}(t)) \hat{V}(t) + b(\hat{m}(t))b(\hat{m}(t))^T \right).$$

starting from $\hat{m}(t_{k-1}) = m_{k-1}$ and $\hat{V}(t_{k-1}) = V_{k-1}$, and we solve the ODE at t_k to get $\hat{m}(t_k) = \bar{m}_k$ and $\hat{V}(t_k) = \bar{V}_k$.

Solve the ODEs with e.g., Euler, Runge-Kutta.

This gives a continuous-discrete extended Kalman filter, which consists in solving a system of non-linear ODEs in the prediction step.

We can also solve ① in a different flavour. Recall:

$$\int X_k P(X_k | Y_{1:k-1}) dX_k$$

$$= \int X_k P(X_k | X_{k-1}) P(X_{k-1} | Y_{1:k-1}) dX_k dX_{k-1}$$

$$m_k^- = \mathbb{E} \left[X(t_k) \mid X(t_{k-1}) = m_{k-1} \right] N(X_k | m_{k-1}, V_{k-1}) dX_{k-1}$$

a Non-linear integrand.

$$\int (X_k - m_k^-)(X_k - m_k^-)^T P(X_k | Y_{1:k-1}) dX_k$$

$$V_k^- \stackrel{\text{def}}{=} \int \mathbb{E}[(X(t_k) - m_k^-)(X(t_k) - m_k^-)^T \mid X(t_{k-1}) = m_{k-1}] N(X_k | m_{k-1}, V_{k-1}) dX_{k-1}$$

~~$\stackrel{\text{def}}{=}$ Cov $\{X(t_k) \mid X(t_{k-1}) = m_{k-1}\}$.~~

Denote $f(x) := \mathbb{E}[X(t_k) \mid X(t_{k-1}) = x]$, then we can linearise f at m_{k-1} by $f(x) \approx f(m_{k-1}) + \boxed{Jf}(m_{k-1})(x - m_{k-1}) + \dots$. Then Jacobian of f .

we have

$$m_k^- = f(m_{k-1}) = \mathbb{E}[X(t_k) \mid X(t_{k-1}) = m_{k-1}],$$

$$V_k^- = \boxed{Jf}(m_{k-1}) V_{k-1} \boxed{Jf}(m_{k-1})^T + \text{Cov}[X(t_k) \mid X(t_{k-1}) = m_{k-1}].$$

How V_k^- is arrived? Law of total covariance

This gives another version of the extended Kalman filter.

Similarly, we can also apply the linearisation to solve for M_k , S_k , and K_k in ②. We linearise $h(x) \approx h(\bar{m}_k) + \boxed{Jh}(\bar{m}_k)$ $(x - \bar{m}_k)$. Then we have

$$M_k = \int h(\bar{x}_k) N(\bar{x}_k | \bar{m}_k, V_k^-) d\bar{x}_k^- \\ \approx h(\bar{m}_k)$$

$$S_k = Jh(\bar{m}_k) V_k^- Jh(\bar{m}_k)^T + \hat{\Sigma}_k$$

$$K_k = V_k^- Jh(\bar{m}_k)^T S_k^{-1}$$

We can summarise the extended Kalman filter in the following algorithm.

Algorithm 38 Extended Kalman filter.

Input: Measurements $y_{1:T}$, initial M_0 and V_0 .

Output: $\{M_k, V_k\}_{k=1}^T$

1. Denote $f(x) := \mathbb{E}[X(t_k) | X(t_{k-1}) = x]$.

1. For $k = 1, 2, \dots, T$ do

$$\bar{m}_k = \mathbb{E}[X(t_k) | X(t_{k-1}) = M_{k-1}]$$

$$F_k = Jf(\bar{m}_{k-1})$$

Denote $f(x) := \mathbb{E}[X(t_k) | X(t_{k-1}) = x]$

$$\Sigma_k = \text{Cov}[X(t_k) | X(t_{k-1}) = m_{k-1}].$$

$$V_k^- = F_k V_{k-1} F_k^T + \Sigma_k.$$

$$H_k = Jh(m_k^-)$$

$$S_k = H_k V_k^- H_k^T + \Sigma_k$$

$$K_k = V_k^- H_k^T S_k^{-1}$$

$$m_k = m_k^- + K_k (y_k - h(m_k^-))$$

$$V_k = V_k^- - K_k S_k K_k^T$$

Return $\{m_k, V_k\}_{k=1}^T$.

Recall Lectures 3 and 4 for how to compute the conditional mean and covariance. E.g., Euler-Maryama.

Mutatis Mutandis, we can derive the extended Kalman smoother

Algorithm 39. Extended Kalman smoother

Input: $y_{1:T}, m_0, V_0$

Output: $\{m_k^s, V_k^s\}_{k=1}^T$

1. Do the extended Kalman filter to obtain

$\{m_k, V_k\}_{k=1}^T$ and $\{m_k^-, V_k^-\}_{k=2}^T$

2. Let ~~m_T^s~~ $m_T^s = m_T$ and $V_T^s = V_T$

3. For $k = T-1, T-2, \dots, 1$ do

$$F_{k+1} = Jf(m_k)$$

$$G_k = V_k F_{k+1}^T (V_{k+1}^-)^{-1}$$

$$M_k^S = M_k + G_k (M_{k+1}^S - M_{k+1}^-)$$

$$V_k^S = V_k + G_k (V_{k+1}^S - V_{k+1}^-) G_k^T$$

return $\{M_k^S, V_k^S\}_{k=1}^T$

Example 40.

$$d \begin{bmatrix} X_1(t) \\ X_2(t) \\ X_3(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ -\left(\frac{\sqrt{3}}{g(X_3(t))}\right)^2 & \frac{-2\sqrt{3}}{g(X_3(t))} & 0 \\ 0 & 0 & -1 \end{bmatrix} \begin{bmatrix} X_1(t) \\ X_2(t) \\ X_3(t) \end{bmatrix} dt + \begin{bmatrix} 0 \\ 2\left(\frac{\sqrt{3}}{g(X_3(t))}\right)^{\frac{3}{2}} \\ 0 \end{bmatrix} dW(t)$$

$$g(x) := \exp(x),$$

$$Y_k = [1 \ 0 \ 0] X(t_k) + \varepsilon_k, \quad \varepsilon_k \sim N(0, I),$$

The non-linear SDE above is a non-stationary extension of the a Matérn Gaussian process. It is useful in modelling non-stationary latent processes.

~~Now we apply the extended~~

Now generate a pair of $(X_{1:T}, Y_{1:T})$ then estimate $X_{1:T}$ by the extended Kalman filter and smoother. We need to take the Jacobian of the transition $\mathbb{E}[X(t_k) | X(t_{k-1}) = x]$.

See 'lec6_extended_kf.ipynb' for how to do so.

Recall the general formulations of Gaussian filtering and smoothing in Algorithms 36 and 37, respectively. We see that ~~these filters and smoothers~~, the essence of them is really solving the Gaussian expectations. For instance, the idea of the extended Kalman filter is to linearise the integrand, so that we can solve the expectations in closed-form.

We can also solve the Gaussian expectations by quadrature, viz.,

$$\int g(x) N(x|m, v) dx \approx \sum_{i=1}^N w_i g(x_i)$$

by a bunch of quadrature nodes and weights $\{x_i, w_i\}_{i=1}^N$. For example, we can let x_i 's be the samples drawn from $N(x|m, v)$. We can also choose them in terms of some orthonormal polynomials, for instance, Gauss-Hermite, unscented transform, and spherical cubature. This gives the celebrated sigma-points filters and smoothers.

Gauss-Hermite quadrature (n -order) is exact for polynomial integrands of degree $\leq 2n-1$.