

BAYESIAN STATE ANALYSIS ON LINEAR GAUSSIAN DYNAMICAL SYSTEMS

Outline:

- ⊕ Bayesian linear regression
- ⊕ Bayesian state analysis on linear Gaussian dynamical systems
 - ◆ Kalman filter
 - ◆ RTS (Rauch-Tung-Striebel) smoother
 - ◆ Backward sampler
- ⊕ Analytical approximation on Bayesian state analysis of nonlinear dynamical systems
 - ◆ Uncertainty propagation
 - ◆ Extended Kalman filter
 - ◆ Unscented Kalman filter

⊕ Bayesian linear regression

Introduction:

Consider the following model:

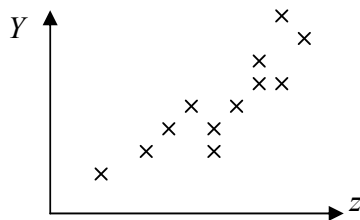
$$Y = aX + E$$

where $E \sim N(0, \Sigma)$, $X \sim N(\mu_X, \Sigma_X)$ and $E \perp X$. We observe $Y = \hat{Y}$ (data), and

we'd like to know $f(x|\hat{Y})$. It is called linear regression because the problem is

linear in the uncertain variable X . Also, note that all uncertainties are Gaussian.

Although this problem is linear, but it can actually handle nonlinearity, e.g.:



$$z_i \in R, \hat{Y}_i \in R, i = 1 \dots N$$

$$g_X(z) = X_0 + X_1 z + X_2 z^2$$

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_N \end{bmatrix} = \begin{bmatrix} 1 & z_1 & z_1^2 \\ \vdots & \vdots & \vdots \\ 1 & z_N & z_N^2 \end{bmatrix} \begin{bmatrix} X_0 \\ X_1 \\ X_2 \end{bmatrix} + \begin{bmatrix} E_1 \\ \vdots \\ E_N \end{bmatrix} \Rightarrow Y = aX + E$$

It turns out that $f(x|\hat{Y})$ is also Gaussian, i.e. $N(\mu_x, \Sigma_x)$ is a conjugate prior. To see this, observe that

$$f(x|\hat{Y}) = \frac{f(\hat{Y}|x)f(x)}{f(\hat{Y})} \propto e^{-\frac{1}{2}(\hat{Y}-ax)^T \Sigma^{-1}(\hat{Y}-ax)} e^{-\frac{1}{2}(x-\mu_x)^T \Sigma_x^{-1}(x-\mu_x)}$$

is log-quadratic in X , so $f(x|\hat{Y})$ is Gaussian. Therefore,

$$f(x|\hat{Y}) = N(\mu_{x|\hat{Y}}, \Sigma_{x|\hat{Y}}), \text{ where}$$

$$\begin{aligned} \mu_{x|\hat{Y}} &= \mu_x + \text{Cov}(X, Y) \text{Cov}(Y)^{-1} (\hat{Y} - E(Y)) = \mu_x + \Sigma_x a^T (a \Sigma_x a^T + \Sigma)^{-1} (\hat{Y} - a \mu_x) \\ \Sigma_{x|\hat{Y}} &= \Sigma_x - \text{Cov}(X, Y) \text{Cov}(Y)^{-1} \text{Cov}(Y, X) = \Sigma_x - \Sigma_x a^T (\Sigma + a \Sigma_x a^T)^{-1} a \Sigma_x \end{aligned}$$

Proof:

Note that at the mean value of the Gaussian PDF $f(x|\hat{Y})$, $\nabla_x f(x|\hat{Y}) = 0$. That

means the mean of the Gaussian PDF can be found by solving $\nabla_x f(x|\hat{Y}) = 0$. Also,

the covariance matrix of the Gaussian PDF $f(x|\hat{Y})$ is equal to

$(-\nabla_x^2 \log[f(x|\hat{Y})])^{-1}$. Therefore, we have

$$\begin{aligned} \mu_{x|\hat{Y}} &= (a^T \Sigma^{-1} a + \Sigma_x^{-1})^{-1} (\Sigma_x^{-1} \mu_x + a^T \Sigma^{-1} \hat{Y}) \\ &= (\Sigma_x a^T \Sigma^{-1} a + I)^{-1} \Sigma_x (\Sigma_x^{-1} \mu_x + a^T \Sigma^{-1} \hat{Y}) \\ &= (\Sigma_x a^T \Sigma^{-1} a + I)^{-1} (\mu_x + \Sigma_x a^T \Sigma^{-1} \hat{Y}) \\ &= (\Sigma_x a^T \Sigma^{-1} a + I)^{-1} (\mu_x + \Sigma_x a^T \Sigma^{-1} \hat{Y} - \Sigma_x a^T \Sigma^{-1} a \mu_x + \Sigma_x a^T \Sigma^{-1} a \mu_x) \end{aligned}$$

$$\begin{aligned}
 &= (\Sigma_X a^T \Sigma^{-1} a + I)^{-1} \left[(\Sigma_X a^T \Sigma^{-1} a + I) \mu_X + \Sigma_X a^T \Sigma^{-1} (\hat{Y} - a \mu_X) \right] \\
 &= \mu_X + (\Sigma_X a^T \Sigma^{-1} a + I)^{-1} \Sigma_X a^T \Sigma^{-1} (\hat{Y} - a \mu_X) \\
 &= \mu_X + (\Sigma_X a^T \Sigma^{-1} a + I)^{-1} \Sigma_X a^T \Sigma^{-1} (\hat{Y} - a \mu_X) \\
 &\because (PQ + I)^{-1} P = P(QP + I)^{-1} \\
 &= \mu_X + \Sigma_X a^T (a \Sigma_X a^T + \Sigma)^{-1} (\hat{Y} - a \mu_X) \\
 \\
 \Sigma_{x|\hat{y}} &= \left(-\nabla_x^2 \log [f(x|\hat{Y})] \right)^{-1} = (a^T \Sigma^{-1} a + \Sigma_X^{-1})^{-1} \\
 &\because (VC^{-1}V^T + A^{-1})^{-1} = A - AV(C + V^T AV)^{-1} V^T A \\
 &= \Sigma_X - \Sigma_X a^T (\Sigma + a \Sigma_X a^T)^{-1} a \Sigma_X
 \end{aligned}$$

✦ Bayesian state analysis on linear Gaussian dynamical systems

Introduction:

For linear dynamical systems with Gaussian uncertainties, the Bayesian state PDF updating can be done analytically. This is because the updated state PDF is Gaussian due to the conjugate priors. Moreover, the state PDF can be updated sequentially. The resulting Bayesian state estimation algorithm is called the Kalman filter.

Consider the following discrete-time linear state-space dynamical system:

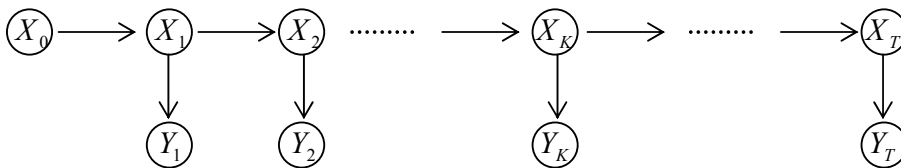
$$X_k = a_{k-1} X_{k-1} + b_{k-1} u_{k-1} + W_{k-1} \quad X_0 \sim N(\mu_{0|0}, \Sigma_{0|0}) \quad (\text{state equation})$$

$$Y_k = c_k X_k + d_k u_k + V_k \quad k = 1, \dots, T \quad (\text{observation equation})$$

where u_k = known system input at time k , W_k, V_k = modeling error, $W_k \sim N(0, \Sigma_W)$,

$V_k \sim N(0, \Sigma_V)$, $W_i \perp W_j (i \neq j)$, $V_i \perp V_j (i \neq j)$, $W_i \perp V_j$, $X_0 \perp W_j$, $X_0 \perp V_j$; a_k ,

b_k , c_k , d_k , Σ_W , Σ_V , $\mu_{0|0}, \Sigma_{0|0}$ are known vectors and matrices. This model class creates a hidden Markov chain.



Note that under this setting, all uncertain variables are jointly Gaussian; also, the model class is linear in all uncertain variables. This implies that if we observe data $\hat{Y} = \{\hat{Y}_1, \dots, \hat{Y}_T\}$, the posterior PDF $f(x | \hat{Y}) \equiv f(x_0, \dots, x_T | \hat{Y})$ is jointly Gaussian, where

$$\begin{aligned} f(x_0, \dots, x_T | \hat{Y}) &= \frac{f(\hat{Y} | x_0, \dots, x_T) f(x_0, \dots, x_T)}{f(\hat{Y})} = \frac{\prod_{k=1}^T f(\hat{Y}_k | x_k) \cdot \prod_{k=1}^T f(x_k | x_{k-1})}{f(\hat{Y})} \\ &= \text{const} \cdot \prod_{k=1}^T e^{-\frac{1}{2}(\hat{Y}_k - c_k x_k - d_k u_k)^T \Sigma_V^{-1} (\hat{Y}_k - c_k x_k - d_k u_k)} \cdot \prod_{k=1}^T e^{-\frac{1}{2}(x_k - a_{k-1} x_{k-1} - b_{k-1} u_{k-1})^T \Sigma_W^{-1} (x_k - a_{k-1} x_{k-1} - b_{k-1} u_{k-1})} \\ &\log \left[f(x_0, \dots, x_T | \hat{Y}) \right] \\ &= \text{const} - \frac{1}{2} \sum_{k=1}^T \left[(\hat{Y}_k - c_k x_k - d_k u_k)^T \Sigma_V^{-1} (\hat{Y}_k - c_k x_k - d_k u_k) \right] \\ &\quad - \frac{1}{2} \sum_{k=1}^T \left[(x_k - a_{k-1} x_{k-1} - b_{k-1} u_{k-1})^T \Sigma_W^{-1} (x_k - a_{k-1} x_{k-1} - b_{k-1} u_{k-1}) \right] \end{aligned}$$

But how do we obtain the mean and the covariance matrix of $f(x | \hat{Y})$? The mean may be simply, i.e. differentiate $\log \left[f(x_0, \dots, x_T | \hat{Y}) \right]$ w.r.t. x_k and solve for zero, we get $E(X_k | \hat{Y})$. But how about the covariance matrix or even $\text{Cov}(X_k | \hat{Y})$? Get the Hessian of $-\log \left[f(x_0, \dots, x_T | \hat{Y}) \right]$ and calculate the inverse? No, we don't want to do so since the inversion will be on a huge matrix.

Another class of interesting problems is to obtain $f(x_k | \hat{Y}_{1:k})$, where $\hat{Y}_{1:k} = \{\hat{Y}_1, \dots, \hat{Y}_k\}$. One can see that this posterior PDF is also Gaussian, where the mean and covariance matrix can, again, be obtained by solving the gradient of $\log \left[f(x_0, \dots, x_k | \hat{Y}_{1:k}) \right]$ for zero and also by calculating the inverse of the Hessian of $-\log \left[f(x_0, \dots, x_k | \hat{Y}_{1:k}) \right]$. But again, we are required to do an inversion on a huge matrix.

Terminology:

A Bayesian filtering problem is to obtain $f(x_k | \hat{Y}_{1:k})$ for all k , while a Bayesian

smoothing problem is to obtain $f(x_k | \hat{Y})$ for all k .

Kalman filter:

Kalman filter provides a smart way of calculating the mean and covariance matrix of $f(x_k | \hat{Y}_{1:k})$ without inverting huge matrices. Basically, Kalman filter is an algorithm that derives $f(x_{k+1} | \hat{Y}_{1:k+1})$ based on the prior $f(x_k | \hat{Y}_{1:k})$ and the new data \hat{Y}_{k+1} , or equivalently, derives $E(X_{k+1} | \hat{Y}_{1:k+1})$ and $Cov(X_{k+1} | \hat{Y}_{1:k+1})$ based on $E(X_k | \hat{Y}_{1:k})$, $Cov(X_k | \hat{Y}_{1:k})$ and the new data \hat{Y}_{k+1} . One can see once this algorithm is finished, we can obtain $f(x_k | \hat{Y}_{1:k})$ recursively starting from $f(x_0 | \hat{Y}_{1:0}) \equiv f(x_0) = N(\mu_{0|0}, \Sigma_{0|0})$.

Let us denote $\mu_{p|q} \equiv E(X_p | \hat{Y}_{1:q})$ and $\Sigma_{p|q} \equiv Cov(X_p | \hat{Y}_{1:q})$.

Algorithm: Kalman filter

1. Starting from $\mu_{0|0}$ and $\Sigma_{0|0}$
2.
$$\begin{aligned} \mu_{k+1|k} &= a_k \mu_{k|k} + b_k u_k \\ \Sigma_{k+1|k} &= a_k \Sigma_{k|k} a_k^T + \Sigma_W \end{aligned} \quad (\text{Uncertainty propagation})$$
3.
$$\begin{aligned} &\mu_{k+1|k+1} \\ &= \mu_{k+1|k} + Cov(X_{k+1}, Y_{k+1} | \hat{Y}_{1:k}) Cov(Y_{k+1} | \hat{Y}_{1:k})^{-1} (\hat{Y}_{k+1} - E(Y_{k+1} | \hat{Y}_{1:k})) \\ &= \mu_{k+1|k} + \Sigma_{k+1|k} c_{k+1}^T (c_{k+1} \Sigma_{k+1|k} c_{k+1}^T + \Sigma_V)^{-1} (\hat{Y}_{k+1} - c_{k+1} \mu_{k+1|k} - d_{k+1} u_{k+1}) \quad (\text{Bayesian update}) \\ &\Sigma_{k+1|k+1} \\ &= \Sigma_{k+1|k} - Cov(X_{k+1}, Y_{k+1} | \hat{Y}_{1:k}) Cov(Y_{k+1} | \hat{Y}_{1:k})^{-1} Cov(Y_{k+1}, X_{k+1} | \hat{Y}_{1:k}) \\ &= \Sigma_{k+1|k} - \Sigma_{k+1|k} c_{k+1}^T (c_{k+1} \Sigma_{k+1|k} c_{k+1}^T + \Sigma_V)^{-1} c_{k+1} \Sigma_{k+1|k} \end{aligned}$$

Note:

$$E(Y_{k+1} | \hat{Y}_{1:k}) = E(c_{k+1} X_{k+1} + d_{k+1} u_{k+1} + V_{k+1} | \hat{Y}_{1:k}) = c_{k+1} \mu_{k+1|k} + d_{k+1} u_{k+1}$$

$$Cov(Y_{k+1} | \hat{Y}_{1:k}) = Cov(c_{k+1} X_{k+1} + d_{k+1} u_{k+1} + V_{k+1} | \hat{Y}_{1:k}) = c_{k+1} \Sigma_{k+1|k} c_{k+1}^T + \Sigma_V$$

$$\begin{aligned}
 & \text{Cov}(X_{k+1}, Y_{k+1} | \hat{Y}_{1:k}) \\
 &= E \left\{ \left[X_{k+1} - E(X_{k+1} | \hat{Y}_{1:k}) \right] \left[Y_{k+1} - E(Y_{k+1} | \hat{Y}_{1:k}) \right]^T \middle| \hat{Y}_{1:k} \right\} \\
 &= E \left\{ \left[X_{k+1} - E(X_{k+1} | \hat{Y}_{1:k}) \right] \left[c_{k+1} X_{k+1} + d_{k+1} u_{k+1} + V_{k+1} - c_{k+1} E(X_{k+1} | \hat{Y}_{1:k}) - d_{k+1} u_{k+1} \right]^T \middle| \hat{Y}_{1:k} \right\} \\
 &= E \left\{ \left[X_{k+1} - E(X_{k+1} | \hat{Y}_{1:k}) \right] \left[c_{k+1} X_{k+1} - c_{k+1} E(X_{k+1} | \hat{Y}_{1:k}) \right]^T \middle| \hat{Y}_{1:k} \right\} \\
 &\quad + E \left\{ \left[X_{k+1} - E(X_{k+1} | \hat{Y}_{1:k}) \right] V_{k+1}^T \middle| \hat{Y}_{1:k} \right\} \\
 &= E \left\{ \left[X_{k+1} - E(X_{k+1} | \hat{Y}_{1:k}) \right] \left[X_{k+1} - E(X_{k+1} | \hat{Y}_{1:k}) \right]^T \middle| \hat{Y}_{1:k} \right\} c_{k+1}^T = \Sigma_{k+1|k} c_{k+1}^T
 \end{aligned}$$

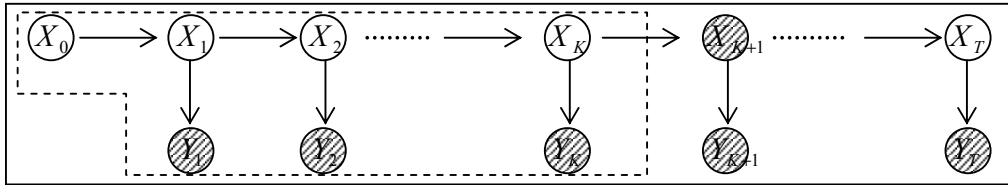
RTS smoother

RTS smoother provides a smart way of calculating the mean and covariance matrix of $f(x_k | \hat{Y})$ (recall that $\hat{Y} = \hat{Y}_{1:T}$) without inverting huge matrices. Basically, RTS smoother starts from the results from Kalman filter and operates backwards in time, i.e. obtain $\mu_{k|T}$ and $\Sigma_{k|T}$ based on $\mu_{k+1|T}$ and $\Sigma_{k+1|T}$. One can see once this algorithm is finished, we can obtain $\mu_{k|T}$ and $\Sigma_{k|T}$ for all k recursively starting from $f(x_T | \hat{Y}) = N(\mu_{T|T}, \Sigma_{T|T})$. Note that $\mu_{T|T}$ and $\Sigma_{T|T}$ are obtained from Kalman filter.

Now we establish the backward recursive equation that relates $\mu_{k|T}$ and $\Sigma_{k|T}$ to

$\mu_{k+1|T}$ and $\Sigma_{k+1|T}$. First note that

$$f(x_k | X_{k+1}, \hat{Y}_{1:T}) = f(x_k | X_{k+1}, \hat{Y}_{1:k}) \equiv N(\mu_{k|k}^*(X_{k+1}), \Sigma_{k|k}^*(X_{k+1}))$$



Moreover,

$$\begin{aligned}
 \boldsymbol{\mu}_{k|k}^* (X_{k+1}) &= \boldsymbol{\mu}_{k|k} + \text{Cov}(X_k, X_{k+1} | \hat{Y}_{1:k}) \text{Cov}(X_{k+1} | \hat{Y}_{1:k})^{-1} (X_{k+1} - E(X_{k+1} | \hat{Y}_{1:k})) \\
 &= \boldsymbol{\mu}_{k|k} + \boldsymbol{\Sigma}_{k|k} \boldsymbol{a}_k^T (\boldsymbol{\Sigma}_{k+1|k})^{-1} (X_{k+1} - \boldsymbol{\mu}_{k+1|k}) \\
 \boldsymbol{\Sigma}_{k|k}^* (X_{k+1}) &= \boldsymbol{\Sigma}_{k|k} - \text{Cov}(X_k, X_{k+1} | \hat{Y}_{1:k}) \text{Cov}(X_{k+1} | \hat{Y}_{1:k})^{-1} \text{Cov}(X_{k+1}, X_k | \hat{Y}_{1:k}) \\
 &= \boldsymbol{\Sigma}_{k|k} - \boldsymbol{\Sigma}_{k|k} \boldsymbol{a}_k^T (\boldsymbol{\Sigma}_{k+1|k})^{-1} \boldsymbol{a}_k \boldsymbol{\Sigma}_{k|k} = \boldsymbol{\Sigma}_{k|k}^*
 \end{aligned}$$

Note:

$$\begin{aligned}
 &\text{Cov}(X_k, X_{k+1} | \hat{Y}_{1:k}) \\
 &= E \left\{ \left[X_k - E(X_k | \hat{Y}_{1:k}) \right] \left[X_{k+1} - E(X_{k+1} | \hat{Y}_{1:k}) \right]^T \middle| \hat{Y}_{1:k} \right\} \\
 &= E \left\{ \left[X_k - E(X_k | \hat{Y}_{1:k}) \right] \left[a_k X_k + b_k u_k + W_k - a_k E(X_k | \hat{Y}_{1:k}) - b_k u_k \right]^T \middle| \hat{Y}_{1:k} \right\} \\
 &= E \left\{ \left[X_k - E(X_k | \hat{Y}_{1:k}) \right] \left[a_k X_k - a_k E(X_k | \hat{Y}_{1:k}) \right]^T \middle| \hat{Y}_{1:k} \right\} + E \left\{ \left[X_k - E(X_k | \hat{Y}_{1:k}) \right] W_k^T \middle| \hat{Y}_{1:k} \right\} \\
 &= E \left\{ \left[X_k - E(X_k | \hat{Y}_{1:k}) \right] \left[X_k - E(X_k | \hat{Y}_{1:k}) \right]^T \middle| \hat{Y}_{1:k} \right\} \boldsymbol{a}_k^T = \boldsymbol{\Sigma}_{k|k} \boldsymbol{a}_k^T
 \end{aligned}$$

Implementing the following identity:

$$E_Y (E_X (X | Y)) = E_X (X) \quad \text{Cov}_X (X) = E_Y [\text{Cov}_X (X | Y)] + \text{Cov}_Y [E_X (X | Y)]$$

One can see that

$$\begin{aligned}
 \boldsymbol{\mu}_{k|T} &= E(X_k | \hat{Y}_{1:T}) = E \left[E(X_k | X_{k+1}, \hat{Y}_{1:T}) \middle| \hat{Y}_{1:T} \right] = E \left[E(X_k | X_{k+1}, \hat{Y}_{1:k}) \middle| \hat{Y}_{1:T} \right] \\
 &= E(\boldsymbol{\mu}_{k|k}^* (X_{k+1}) | \hat{Y}_{1:T}) = E \left(\boldsymbol{\mu}_{k|k} + \boldsymbol{\Sigma}_{k|k} \boldsymbol{a}_k^T (\boldsymbol{\Sigma}_{k+1|k})^{-1} (X_{k+1} - \boldsymbol{\mu}_{k+1|k}) \middle| \hat{Y}_{1:T} \right) \\
 &= \boldsymbol{\mu}_{k|k} + \boldsymbol{\Sigma}_{k|k} \boldsymbol{a}_k^T (\boldsymbol{\Sigma}_{k+1|k})^{-1} (\boldsymbol{\mu}_{k+1|T} - \boldsymbol{\mu}_{k+1|k}) \\
 \boldsymbol{\Sigma}_{k|T} &= \text{Cov}(X_k | \hat{Y}_{1:T}) = E \left[\text{Cov}(X_k | X_{k+1}, \hat{Y}_{1:T}) \middle| \hat{Y}_{1:T} \right] + \text{Cov} \left[E(X_k | X_{k+1}, \hat{Y}_{1:T}) \middle| \hat{Y}_{1:T} \right] \\
 &= E \left[\text{Cov}(X_k | X_{k+1}, \hat{Y}_{1:k}) \middle| \hat{Y}_{1:T} \right] + \text{Cov} \left[E(X_k | X_{k+1}, \hat{Y}_{1:k}) \middle| \hat{Y}_{1:T} \right] \\
 &= \boldsymbol{\Sigma}_{k|k}^* + \text{Cov} \left[\boldsymbol{\mu}_{k|k}^* (X_{k+1}) \middle| \hat{Y}_{1:T} \right] \\
 &= \boldsymbol{\Sigma}_{k|k} - \boldsymbol{\Sigma}_{k|k} \boldsymbol{a}_k^T (\boldsymbol{\Sigma}_{k+1|k})^{-1} \boldsymbol{a}_k \boldsymbol{\Sigma}_{k|k} + \boldsymbol{\Sigma}_{k|k} \boldsymbol{a}_k^T (\boldsymbol{\Sigma}_{k+1|k})^{-1} \boldsymbol{\Sigma}_{k+1|T} (\boldsymbol{\Sigma}_{k+1|k})^{-1} \boldsymbol{a}_k \boldsymbol{\Sigma}_{k|k}
 \end{aligned}$$

Note that the data is not needed for the RTS smoother.

Algorithm: RTS smoother

1. Run Kalman filter first
2. Starting from $\boldsymbol{\mu}_{T|T}$ and $\boldsymbol{\Sigma}_{T|T}$
3. Operate backwards in time

$$\begin{aligned}\mu_{k|T} &= \mu_{k|k} + \Sigma_{k|k} a_k^T (\Sigma_{k+1|k})^{-1} (\mu_{k+1|T} - \mu_{k+1|k}) \\ \Sigma_{k|T} &= \Sigma_{k|k} - \Sigma_{k|k} a_k^T (\Sigma_{k+1|k})^{-1} a_k \Sigma_{k|k} + \Sigma_{k|k} a_k^T (\Sigma_{k+1|k})^{-1} \Sigma_{k+1|T} (\Sigma_{k+1|k})^{-1} a_k \Sigma_{k|k}\end{aligned}$$

Backward sampler:

In many times, we are interested in the maximum state response over the entire time interval $[0, T]$. However, from the results of Kalman filter and RTS smoother, we lose the correlation information between the states of different time. This correlation information is essential for understanding the maximum state response over time. We describe an algorithm that draws state time history samples from the posterior PDF $f(x_0, \dots, x_T | \hat{Y})$. This algorithm requires us to run Kalman filter first.

Note that

$$\begin{aligned}& f(x_0, \dots, x_T | \hat{Y}) \\ &= f(x_T | \hat{Y}) f(x_{T-1} | x_T, \hat{Y}) \cdots f(x_K | x_{K+1}, \dots, x_T, \hat{Y}) \cdots f(x_0 | x_1, \dots, x_T, \hat{Y}) \\ &= f(x_T | \hat{Y}_{1:T}) f(x_{T-1} | x_T, \hat{Y}_{1:T-1}) \cdots f(x_k | x_{k+1}, \hat{Y}_{1:k}) \cdots f(x_1 | x_2, \hat{Y}_1) f(x_0 | x_1)\end{aligned}$$

A strategy is to first sample \hat{X}_T from $f(x_T | \hat{Y}) = N(\mu_{T|T}, \Sigma_{T|T})$, then sample \hat{X}_{T-1} from $f(x_{T-1} | \hat{X}_T, \hat{Y}_{1:T-1}) = N(\mu_{T-1|T-1}^*(\hat{X}_T), \Sigma_{T-1|T-1}^*)$, then sample \hat{X}_{T-2} from $N(\mu_{T-2|T-2}^*(\hat{X}_{T-1}), \Sigma_{T-2|T-2}^*)$ and so on to get a sample of the state time history. Do this many times independently to get independent state time history samples. Afterwards, we can use these samples to estimate the expected value of the maximum state response based on the Law of Large Number. Note that the data is not needed for the backward sampler.

Algorithm: Backward sampler

1. Do Kalman filter first

$$\hat{X}_T \sim N(\mu_{T|T}, \Sigma_{T|T})$$

2. Sample $\hat{X}_{T-1} \sim N\left(\begin{array}{c} \mu_{T-1|T-1} + \Sigma_{T-1|T-1} a_{T-1}^T (\Sigma_{T|T-1})^{-1} \cdot (\hat{X}_T - x_{T|T-1}), \\ \Sigma_{T-1|T-1} + \Sigma_{T-1|T-1} a_{T-1}^T (\Sigma_{T|T-1})^{-1} a_{T-1} \Sigma_{T-1|T-1} \end{array}\right)$

⋮

3. Do (2) N times to get N i.i.d. samples from $f(x_0, \dots, x_T | \hat{Y})$

⊕ Analytical approximation on Bayesian state analysis of nonlinear dynamical systems

Introduction:

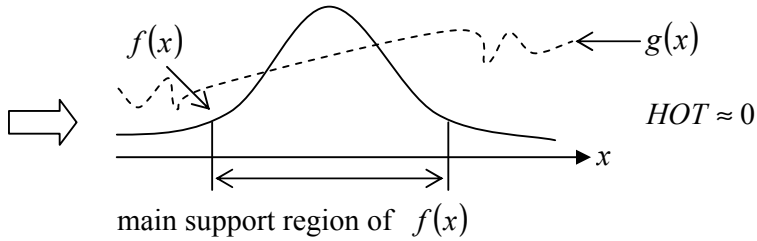
Many dynamical models are nonlinear or non-Gaussian. In this case, the above analysis breaks down. However, we can linearize and Gaussianize those models so Kalman filter, RTS smoother and backward sampler can still work approximately.

Extended Kalman filter: linearization on uncertainty propagation for $Y = g(X)$

Do Taylor expansion on $g(X)$ around $E(X)$, we get

$$Y = g(X) = g(E(X)) + \nabla_x g(E(X)) \cdot (X - E(X)) + HOT$$

Assuming $g(\cdot)$ is roughly linear in the main support region of the PDF $f(x)$,



we have

$$Y = g(X) \approx g(E(X)) + \nabla_x g(E(X)) \cdot (X - E(X))$$

$$E(Y) \approx g(E(X)) \text{ and } Cov(Y) \approx \nabla_x g \Big|_{x=E(X)} \cdot Cov(X) \cdot \nabla_x g \Big|_{x=E(X)}^T$$

Note that the truncation error is of 2nd order, and the approximation will be poor if $g(\cdot)$ is highly nonlinear in the main support region of the PDF $f(x)$. This method of approximately propagating the first two moments is sometimes called the First-Order-Second-Moment (FOSM) method. Under this approximation, any nonlinear models can be linearized. The resulting Bayesian filtering algorithm is called extended Kalman filter. We'll skip the smoothing and backward sampling part.

For non-linear models, we can linearize the state and observation equations:

$$\begin{aligned} X_{k+1} &= \phi_k(X_k, u_k, W_k) \\ &\approx \phi_k(\mu_{k|k}^{EKF}, u_k, 0) + \nabla_{X_k} \phi_k(\mu_{k|k}^{EKF}, u_k, 0) \cdot (X_k - \mu_{k|k}^{EKF}) + \nabla_{W_k} \phi_k(\mu_{k|k}^{EKF}, u_k, 0) \cdot W_k \\ &= \underbrace{\nabla_{X_k} \phi_k(\mu_{k|k}^{EKF}, u_k, 0)}_{a_k} \cdot X_k + \underbrace{\left[\phi_k(\mu_{k|k}^{EKF}, u_k, 0) - \nabla_{X_k} \phi_k(\mu_{k|k}^{EKF}, u_k, 0) \cdot \mu_{k|k}^{EKF} \right]}_{b_k u_k} + W_k^{EKF} \end{aligned}$$

$$\begin{aligned}
 Y_k &= \varphi_k(X_k, u_k, V_k) \\
 &\approx \varphi_k(\mu_{k|k-1}^{EKF}, u_k, 0) + \nabla_{X_k} \varphi_k(\mu_{k|k-1}^{EKF}, u_k, 0) \cdot (X_k - \mu_{k|k-1}^{EKF}) + \nabla_{V_k} \varphi_k(\mu_{k|k-1}^{EKF}, u_k, 0) \cdot V_k \\
 &= \underbrace{\nabla_{X_k} \varphi_k(\mu_{k|k-1}^{EKF}, u_k, 0)}_{c_k} \cdot X_k + \underbrace{\left[\varphi_k(\mu_{k|k-1}^{EKF}, u_k, 0) - \nabla_{X_k} \varphi_k(\mu_{k|k-1}^{EKF}, u_k, 0) \cdot \mu_{k|k-1}^{EKF} \right]}_{d_k u_k} + V_k^{EKF}
 \end{aligned}$$

Then proceed with the Kalman filter algorithm. The resulting algorithm is called the extended Kalman filter.

Unscented Kalman filter: propagate moments by matching moments

Let $s = (X - EX) / \|X - EX\|$, consider the Taylor series expansion in the s direction:

$$Y = g(X) = g(E(X)) + \sum_{i=1}^{\infty} \frac{1}{i!} \left. \frac{\partial^i g(x)}{\partial s^i} \right|_{x=E(X)} \cdot \|X - E(X)\|^i$$

where

$$\begin{aligned}
 \left. \frac{\partial^i g(x)}{\partial s^i} \right|_{x=E(X)} &= (s \cdot \nabla_x)^i g(x) \Big|_{x=E(X)} \\
 &= \left(s_1 \frac{\partial}{\partial x_1} + \dots + s_n \frac{\partial}{\partial x_n} \right)^i g(x) \Big|_{x=E(X)} = \left(\sum_{j=1}^n \frac{X_j - E(X_j)}{\|X - EX\|} \cdot (\partial / \partial x_j) \right)^i g(x) \Big|_{x=E(X)}
 \end{aligned}$$

So we get

$$\begin{aligned}
 Y &= g(X) \\
 &= g(E(X)) + \sum_{i=1}^{\infty} \frac{1}{i!} \left(\sum_{j=1}^n \frac{X_j - E(X_j)}{\|X - EX\|} \cdot (\partial / \partial x_j) \right)^i \cdot g(E(X)) \cdot \|X - E(X)\|^i \\
 &= g(E(X)) + \sum_{i=1}^{\infty} \frac{1}{i!} \left(\sum_{j=1}^n [X_j - E(X_j)] \cdot (\partial / \partial x_j) \right)^i \cdot g(E(X))
 \end{aligned}$$

Therefore,

$$E(Y) = g(E(X)) + \sum_{i=1}^{\infty} \frac{1}{i!} E \left[\left(\sum_{j=1}^n [X_j - E(X_j)] \cdot (\partial / \partial x_j) \right)^i \right] \cdot g(E(X))$$

If $g(\cdot)$ is an $(2p-1)^{th}$ order polynomial and if we can find another uncertain variable $\bar{X} (\neq X)$ whose first $2p-1$ moments are identical to those of X . Define $\bar{Y} = g(\bar{X})$. It is clear that $E(\bar{Y}) = E(Y)$. This is because when $g(\cdot)$ is an $(2p-1)^{th}$ order polynomial and the first $2p-1$ moment of \bar{X} and X are identical,

$$\begin{aligned}
 E(Y) &= g(E(X)) + \sum_{i=1}^{2p-1} \frac{1}{i!} E \left[\left(\sum_{j=1}^n [X_j - E(X_j)] \cdot (\partial/\partial x_j) \right)^i \right] \cdot g(E(X)) \\
 &= g(E(\bar{X})) + \sum_{i=1}^{2p-1} \frac{1}{i!} E \left[\left(\sum_{j=1}^n [\bar{X}_j - E(\bar{X}_j)] \cdot (\partial/\partial x_j) \right)^i \right] \cdot g(E(\bar{X})) = E(\bar{Y})
 \end{aligned}$$

Note that given the $2p-1$ moments of X , \bar{X} is not unique. A convenient choice is to take \bar{X} with the following PDF:

$$f(\bar{x}) = \sum_{i=1}^p w_i \delta(\bar{x} - \lambda_i), \quad \sum_{i=1}^p w_i = 1$$


where $\{(\lambda_i, w_i) : i=1, \dots, p\}$ are the locations and weights of the delta functions. We can adjust the $2p-1$ free parameters to match the first $2p-1$ moments of X . Now it is clear that

$$E(Y) = E(\bar{Y}) = \sum_{i=1}^p w_i g(\lambda_i)$$

is an exact solution!! We name this approach as the moment matching method (MM).

Consider the case that we would like to estimate the mean and variance of a scalar function $h(X)$, i.e. we want to estimate $E[h(X)]$ and $\text{Var}[h(X)]$, where X is also a scalar. Also consider the usual case that we cannot analytically determine the gradient of $h(X)$. For the linearization approach (FOSM), we usually need to evaluate $h(\cdot)$ function at three points to find $E[h(X)]$ and $\text{Var}[h(X)]$ since

$$E(h(X)) \approx h(E(X)) \quad \text{and} \quad \text{Var}(h(X)) \approx \left(\frac{h(E(X) + \Delta x) - h(E(X) - \Delta x)}{2\Delta x} \right)^2 \cdot \text{Var}(X)$$

The FOSM estimates are exact if the $h(\cdot)$ function is linear in the main support region of the PDF of X .

With the same computation cost, we can employ a three-point moment matching method, i.e. let \bar{X} be the uncertain variable with the following PDF:

$$f(\bar{x}) = \sum_{i=1}^3 w_i \delta(\bar{x} - \lambda_i), \quad \sum_{i=1}^3 w_i = 1$$

As discussed in the above, we can match the first five moments of X using \bar{X} . So we have

$$E(h(X)) \approx \sum_{i=1}^3 w_i h(\lambda_i) \quad \text{and} \quad E(h(X)^2) \approx \sum_{i=1}^3 w_i h(\lambda_i)^2$$

The MM estimates for $E(h(X))$ is exact if the $h(\cdot)$ function is a fifth-order polynomial (or less) in the main support region of the PDF of X ; the $E(h(X)^2)$ estimate is exact if the $h(\cdot)$ function is a quadratic polynomial (or less) in the main support region of the PDF of X .

You may wonder how to match moments by selecting appropriate $\{(\lambda_i, w_i) : i = 1, \dots, p\}$? It is really simple: just solve the following equations:

$$\sum_{i=1}^p w_i = 1 \quad \sum_{i=1}^p w_i \lambda_i = E(X) \quad \dots \quad \sum_{i=1}^p w_i \lambda_i^{2p-1} = E(X^{2p-1})$$

When X is some standard uncertain variable, we don't even need to solve them:

1. If X is uniform, $\{(\lambda_i, w_i) : i = 1, \dots, p\}$ are related to the locations and weights of Gauss-Legendre quadrature. See the following link:
<http://mathworld.wolfram.com/Legendre-GaussQuadrature.html>
2. If X is Gaussian, $\{(\lambda_i, w_i) : i = 1, \dots, p\}$ are related to the locations and weights of Gauss-Hermite quadrature. See the following link:
<http://mathworld.wolfram.com/Hermite-GaussQuadrature.html>
3. If X is exponential, $\{(\lambda_i, w_i) : i = 1, \dots, p\}$ are related to the locations and weights of Gauss-Laguerre quadrature. See the following link:
<http://mathworld.wolfram.com/Laguerre-GaussQuadrature.html>

If we use MM to propagate the first two moments $E(X_k | \hat{Y}_{1:k})$ and $Cov(X_k | \hat{Y}_{1:k})$ for approximate Bayesian filtering, the resulting algorithm is called the unscented Kalman filter [1].

Maximum entropy argument

We can think of the extended/unscented Kalman filters as first-two-moment Bayesian filters based on the maximum entropy principle. In the case that we only want to propagate the first two moments $E(X_k | \hat{Y}_{1:k})$ and $Cov(X_k | \hat{Y}_{1:k})$ for nonlinear

dynamical systems, we can argue that the maximum entropy PDF constrained by the two moments $E(X_k | \hat{Y}_{1:k})$ and $Cov(X_k | \hat{Y}_{1:k})$ is Gaussian.

Reference:

[1] Julier, S.J, and Uhlmann, J.K. (1997) "A new extension of the Kalman filter to nonlinear systems." In *Proceedings of AeroSense: The 11th International Symposium on Aerospace/Defense Sensing, Simulation and Controls*.