

BAYESIAN STATE ANALYSIS ON NONLINEAR DYNAMICAL SYSTEMS

Outline:

- ⊕ Particle filter
- ⊕ Particle smoother

- ⊕ Particle filter

Introduction:

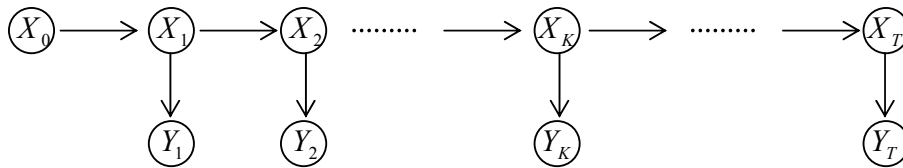
In the case that the state-space model class is nonlinear, Kalman filter and RTS smoother breaks down. Although it is always possible to linearize the nonlinear model so that Kalman filter and RTS smoother can still apply approximately, they can be not reliable. On the other hand, stochastic simulation approaches are not limited to linear model classes and can be adopted to draw samples of the state conditioning on the observation even if the model is nonlinear. One of such stochastic simulation approaches is called particle filter.

Let's now consider a more general version of state-space equations:

$$X_k = \phi_{k-1}(X_{k-1}, u_k, W_k) \quad Y_k = \varphi_k(X_k, u_k, V_k)$$

where u_k = known system input at time k , W_k, V_k = modeling error, $W_i \perp W_j (i \neq j)$,

$V_i \perp V_j (i \neq j)$, $W_i \perp V_j$, $X_0 \perp W_j$, $X_0 \perp V_j$. Note that this model class creates a hidden Markov chain.



For the sake of simplicity, let's consider a simpler model class in this lecture:

$$X_k = \phi_{k-1}(X_{k-1}, u_k) + W_k \quad Y_k = \varphi_k(X_k, u_k) + V_k$$

where $W_k \sim N(0, \Sigma_W)$, $V_k \sim N(0, \Sigma_V)$. Similar to Kalman filter, the goal of particle

filter is to draw samples from $f(x_k | \hat{Y}_{1:k})$, where $\hat{Y}_{1:k} = \{\hat{Y}_1, \dots, \hat{Y}_k\}$ is the data up to time k . Different from Kalman filter, this posterior PDF is not Gaussian. In the

following, we adopt the following notation: $x_{1:k} = \{x_1, \dots, x_k\}$ and $X_{1:k} = \{X_1, \dots, X_k\}$.

The particle filter algorithm provides a way of obtaining samples of $f(x_{1:k+1} | \hat{Y}_{1:k+1})$

given the samples of $f(x_{1:k} | \hat{Y}_{1:k})$ and the new observation \hat{Y}_{k+1} . One can see that

once this algorithm is finished, we can obtain the samples of $f(x_{1:k} | \hat{Y}_{1:k})$ for all k

recursively starting from the samples of $f(x_0 | \hat{Y}_{1:0}) \equiv f(x_0)$. Note that once we get

the samples $\{\hat{X}_{1:k}^i : i = 1, \dots, N\}$ from $f(x_{1:k} | \hat{Y}_{1:k})$, we just keep the $\{\hat{X}_k^i : i = 1, \dots, N\}$

part of the samples; that will give us the samples that are distributed as $f(x_k | \hat{Y}_{1:k})$.

Before we describe the algorithm, let's look at the Bayesian derivations. Given the

samples $\{\hat{X}_{1:k}^i : i = 1, \dots, N\}$, $f(x_{1:k} | \hat{Y}_{1:k})$ can be represented approximately by

$$f(x_{1:k} | \hat{Y}_{1:k}) \approx \frac{1}{N} \sum_{i=1}^N \delta(x_{1:k} - \hat{X}_{1:k}^i)$$

Note that this approximation is in the sense of the Law of Large Number, i.e.

$$E(g(X_{1:k}) | \hat{Y}_{1:k}) \approx \frac{1}{N} \sum_{i=1}^N g(\hat{X}_{1:k}^i) = \int g(x_{1:k}) \cdot \frac{1}{N} \sum_{i=1}^N \delta(x_{1:k} - \hat{X}_{1:k}^i) dx_{1:k}$$

From the Bayes rule, we have

$$\begin{aligned} f(x_{1:k+1} | \hat{Y}_{1:k+1}) &= \frac{f(x_{1:k+1}, \hat{Y}_{1:k+1})}{f(\hat{Y}_{1:k+1})} = \frac{f(x_{1:k}, x_{k+1}, \hat{Y}_{1:k}, \hat{Y}_{k+1})}{f(\hat{Y}_{1:k}, \hat{Y}_{k+1})} \\ &= \frac{f(x_{k+1}, \hat{Y}_{k+1} | x_{1:k}, \hat{Y}_{1:k})}{f(\hat{Y}_{1:k}, \hat{Y}_{k+1})} \cdot f(x_{1:k}, \hat{Y}_{1:k}) \\ &= \frac{f(x_{k+1}, \hat{Y}_{k+1} | x_{1:k}, \hat{Y}_{1:k})}{f(\hat{Y}_{k+1} | \hat{Y}_{1:k})} \cdot f(x_{1:k} | \hat{Y}_{1:k}) \\ &= \frac{f(\hat{Y}_{k+1} | x_{k+1}, x_{1:k}, \hat{Y}_{1:k}) f(x_{k+1} | x_{1:k}, \hat{Y}_{1:k})}{f(\hat{Y}_{k+1} | \hat{Y}_{1:k})} \cdot f(x_{1:k} | \hat{Y}_{1:k}) \\ &= \frac{f(\hat{Y}_{k+1} | x_{k+1}) f(x_{k+1} | x_k)}{f(\hat{Y}_{k+1} | \hat{Y}_{1:k})} \cdot f(x_{1:k} | \hat{Y}_{1:k}) \\ &\approx \frac{f(\hat{Y}_{k+1} | x_{k+1}) f(x_{k+1} | x_k)}{f(\hat{Y}_{k+1} | \hat{Y}_{1:k})} \cdot \frac{1}{N} \sum_{i=1}^N \delta(x_{1:k} - \hat{X}_{1:k}^i) \end{aligned}$$

Note that the new data \hat{Y}_{k+1} plays a major role in the equations. There is in fact a

simple way of drawing samples from $f(x_{1:k+1} | \hat{Y}_{1:k+1})$: (1) draw a sample $\hat{X}_{1:k}$ from

$$f(x_{1:k} | \hat{Y}_{1:k}) \approx \frac{1}{N} \sum_{i=1}^N \delta(x_{1:k} - \hat{X}_{1:k}^i) \quad \text{and} \quad (2) \quad \text{draw a sample } \hat{X}_{k+1} \text{ from}$$

$$\frac{f(\hat{Y}_{k+1} | x_{k+1}) f(x_{k+1} | \hat{X}_k)}{f(\hat{Y}_{k+1} | \hat{Y}_{1:k})} \text{ so that the combined sample } \hat{X}_{1:k+1} \equiv \{\hat{X}_{1:k}, \hat{X}_{k+1}\} \text{ should}$$

be approximately distributed as $f(x_{1:k+1} | \hat{Y}_{1:k+1})$. However, this is not doable directly

because the PDF $\frac{f(\hat{Y}_{k+1} | x_{k+1}) f(x_{k+1} | \hat{X}_k)}{f(\hat{Y}_{k+1} | \hat{Y}_{1:k})}$ cannot be sampled directly.

On the other hand, one can employ the sample-importance resampling (SIR) technique to sample from $f(x_{1:k+1} | \hat{Y}_{1:k+1})$. Let us choose the following importance sampling PDF:

$$q(x_{1:k+1}) = q_0(x_{k+1} | x_k, \hat{Y}_{k+1}) \cdot \frac{1}{N} \sum_{i=1}^N \delta(x_{1:k} - \hat{X}_{1:k}^i)$$

where $q_0(x_{k+1} | x_k, \hat{Y}_{k+1})$ is some PDF of X_{k+1} . The sequential way of doing SIR through time is called the particle filter.

Algorithm: particle filter

1. Let $\{\hat{X}_{1:k}^i : i=1, \dots, N\}$ be the samples distributed as $f(x_{1:k} | \hat{Y}_{1:k})$.
2. Sample $\{\bar{X}_{1:k+1}^i : i=1, \dots, N\}$ from $q(x_{1:k+1})$, i.e. draw $\bar{X}_{1:k}^i$ from $\frac{1}{N} \sum_{i=1}^N \delta(x_{1:k} - \hat{X}_{1:k}^i)$ and draw \bar{X}_{k+1}^i from $q_0(x_{k+1} | \bar{X}_k^i, \hat{Y}_{k+1})$, and do this for $i = 1, \dots, N$.
3. Compute the normalized importance weights:

$$w^{(i)} \equiv \frac{f(\hat{Y}_{k+1} | \bar{X}_{k+1}^i) f(\bar{X}_{k+1}^i | \bar{X}_k^i)}{q_0(\bar{X}_{k+1}^i | \bar{X}_k^i, \hat{Y}_{k+1})} \bigg/ \sum_{j=1}^N \frac{f(\hat{Y}_{k+1} | \bar{X}_{k+1}^j) f(\bar{X}_{k+1}^j | \bar{X}_k^j)}{q_0(\bar{X}_{k+1}^j | \bar{X}_k^j, \hat{Y}_{k+1})}$$

4. Resampling: Let $\hat{X}_{1:k+1}^j = \bar{X}_{1:k+1}^i$ with probability $w^{(i)}$ for $j=1, \dots, N$, then

$\{\hat{X}_{1:k+1}^j : j=1, \dots, N\}$ will be approximately distributed as $f(x_{1:k+1} | \hat{Y}_{1:k+1})$.

Remarks:

1. Since the particle filter algorithm is based on SIR, so it is not good for high dimensional X_k .

2. From the algorithm, it is clear that the samples before time k $\{\hat{X}_{1:k-1}^i : i=1, \dots, N\}$

are not needed, so in the real application, the particle filter algorithm can be a full online algorithm, i.e. no need to store the samples older than the previous time step in the memory. So we can actually change the algorithm into the following:

I. Let $\{\hat{X}_k^i : i=1, \dots, N\}$ be the samples distributed as $f(x_k | \hat{Y}_{1:k})$.

II. Draw \bar{X}_k^i from $\frac{1}{N} \sum_{i=1}^N \delta(x_k - \hat{X}_k^i)$ and draw \bar{X}_{k+1}^i from

$q_0(x_{k+1} | \bar{X}_k^i, \hat{Y}_{k+1})$, and do this for $i=1, \dots, N$.

III. Compute the normalized importance weights:

$$w^{(i)} \equiv \frac{f(\hat{Y}_{k+1} | \bar{X}_{k+1}^i) f(\bar{X}_{k+1}^i | \bar{X}_k^i)}{q_0(\bar{X}_{k+1}^i | \bar{X}_k^i, \hat{Y}_{k+1})} \bigg/ \sum_{j=1}^N \frac{f(\hat{Y}_{k+1} | \bar{X}_{k+1}^j) f(\bar{X}_{k+1}^j | \bar{X}_k^j)}{q_0(\bar{X}_{k+1}^j | \bar{X}_k^j, \hat{Y}_{k+1})}$$

IV. Resampling: Let $\hat{X}_{k+1}^j = \bar{X}_{k+1}^i$ with probability $w^{(i)}$ for $j=1, \dots, N$, then

$\{\hat{X}_{k+1}^j : j=1, \dots, N\}$ will be approximately distributed as $f(x_{k+1} | \hat{Y}_{1:k+1})$.

3. We can replace the random sampling of \bar{X}_k^i in Step II with a deterministic one:

just set $\bar{X}_k^i = \hat{X}_k^i$ for all i .

4. One can, in principle, skip the resampling step, but when doing so, the importance weights need to be accumulated to the next time step. However, this is not feasible because if we do so, the weights will ultimately become highly non-uniform so the effective number of samples will be very small. Therefore, the resampling step is essential. One way of interpreting the resampling step is to consider it as a splitting step according to the sample importance, i.e. the samples with large weights will be duplicated and propagated continuously, but the samples with small weights may be eliminated.

5. In fact, it can be shown that the particle filter algorithm can asymptotically draw

samples that are distributed as $f(x_k | \hat{Y}_{1:k})$. This is because SIR is asymptotically correct.

6. A critical question in the particle filter algorithm is how to select the importance PDF $q_0(x_{k+1} | \bar{X}_k^i, \hat{Y}_{k+1})$. A popular choice is to let $q_0(x_{k+1} | \bar{X}_k^i, \hat{Y}_{k+1})$ be equal to $f(x_{k+1} | \bar{X}_k^i)$ (a Gaussian PDF in our example) because we know how to directly sample and evaluate it. Although this is a convenient choice, but the resulting importance weights may be highly non-uniform. This undesirable situation may occur when the main support region of $f(x_{k+1} | \bar{X}_k^i)$ is very different from that of $f(\hat{Y}_{k+1} | x_{k+1})f(x_{k+1} | \bar{X}_k^i)$.

The optimal choice is to let $q_0(x_{k+1} | \bar{X}_k^i, \hat{Y}_{k+1}) = \frac{f(\hat{Y}_{k+1} | x_{k+1})f(x_{k+1} | \hat{X}_k)}{f(\hat{Y}_{k+1} | \hat{Y}_{1:k})}$

because the importance weights will be completely uniform. However, it is not feasible because we don't know how to directly sample and evaluate it.

An interesting choice is to let $q_0(x_{k+1} | \bar{X}_k^i, \hat{Y}_{k+1})$ be a Gaussian PDF that is

similar to $\frac{f(\hat{Y}_{k+1} | x_{k+1})f(x_{k+1} | \hat{X}_k)}{f(\hat{Y}_{k+1} | \hat{Y}_{1:k})}$. This important sampling PDF can be

obtained using Extend Kalman filter (EKF): Note that the prior PDF $f(x_{k+1} | \hat{X}_k)$ is Gaussian. If we employ the linearization technique in EKF to linearize the likelihood function $f(\hat{Y}_{k+1} | x_{k+1})$ so that it's linear in x_{k+1} , we can obtain a

Gaussian approximation of the posterior PDF $\frac{f(\hat{Y}_{k+1} | x_{k+1})f(x_{k+1} | \hat{X}_k)}{f(\hat{Y}_{k+1} | \hat{Y}_{1:k})}$ and use

that as our importance sampling PDF $q_0(x_{k+1} | \bar{X}_k^i, \hat{Y}_{k+1})$.

7. After we obtain samples $\{\hat{X}_k^i : i = 1, \dots, N\}$, $E[g(X_k) | \hat{Y}_{1:k}]$ can be estimated as

$$E[g(X_k) | \hat{Y}_{1:k}] \approx \frac{1}{N} \sum_{i=1}^N g(\hat{X}_k^i)$$

✦ Particle smoother

Introduction:

The goal of particle smoother is to sample from $f(x_{1:T} | \hat{Y}_{1:T})$, i.e. draw the complete time history of the state conditioning on the data. Its role is very similar to the backward sampler for linear Gaussian model classes, but particle smoother is applicable to nonlinear model classes. Note that if we know how to sample from $f(x_{1:T} | \hat{Y}_{1:T})$, we know how to sample from $f(x_k | \hat{Y}_{1:T})$ for every k , i.e. discard all samples except the state sample at time k . Note that the particle smoother requires the results from particle filter: $\{\hat{X}_k^i : i = 1, \dots, N, k = 1, \dots, T\}$, where \hat{X}_k^i is distributed as

$$f(x_k | \hat{Y}_{1:k}).$$

Consider the following equation:

$$\begin{aligned} & f(x_{1:T} | \hat{Y}_{1:T}) \\ &= f(x_T | \hat{Y}_{1:T}) f(x_{T-1} | x_T, \hat{Y}_{1:T}) \cdots f(x_k | x_{k+1}, x_{k+2}, \dots, x_T, \hat{Y}_{1:T}) \cdots f(x_0 | x_1, \dots, x_T, \hat{Y}_{1:T}) \\ &= f(x_T | \hat{Y}_{1:T}) f(x_{T-1} | x_T, \hat{Y}_{1:T-1}) \cdots f(x_k | x_{k+1}, \hat{Y}_{1:k}) \cdots f(x_0 | x_1) \end{aligned}$$

One way of sampling from $f(x_{1:T} | \hat{Y}_{1:T})$ is to first draw \hat{X}_T^S (the ‘S’ superscript denotes smoother samples to differentiate them from the samples from the particle filter) from $f(x_T | \hat{Y}_{1:T})$ and then draw \hat{X}_{T-1}^S from $f(x_{T-1} | \hat{X}_T^S, \hat{Y}_{1:T-1})$, then draw \hat{X}_{T-2}^S from $f(x_{T-2} | \hat{X}_{T-1}^S, \hat{Y}_{1:T-2})$, and so on to obtain $\hat{X}_{0:T}^S$, distributed as $f(x_{1:T} | \hat{Y}_{1:T})$. Drawing samples from $f(x_{1:T} | \hat{Y}_{1:T})$ is easy with the results from the particle filter; however, $f(x_{T-1} | \hat{X}_T^S, \hat{Y}_{1:T-1})$ cannot be directly sampled.

Consider the following equation:

$$\begin{aligned} f(x_{T-1} | \hat{X}_T^S, \hat{Y}_{1:T-1}) &= \frac{f(x_{T-1}, \hat{X}_T^S, \hat{Y}_{1:T-1})}{f(\hat{X}_T^S, \hat{Y}_{1:T-1})} = \frac{f(\hat{X}_T^S | x_{T-1}, \hat{Y}_{1:T-1}) f(x_{T-1}, \hat{Y}_{1:T-1})}{f(\hat{X}_T^S, \hat{Y}_{1:T-1})} \\ &= \frac{f(\hat{X}_T^S | x_{T-1})}{f(\hat{X}_T^S | \hat{Y}_{1:T-1})} \cdot f(x_{T-1} | \hat{Y}_{1:T-1}) \approx \frac{f(\hat{X}_T^S | x_{T-1})}{f(\hat{X}_T^S | \hat{Y}_{1:T-1})} \cdot \frac{1}{N} \sum_{i=1}^N \delta(x_{T-1} - \hat{X}_{T-1}^i) \end{aligned}$$

$$\begin{aligned}
 &\approx \sum_{i=1}^N \frac{1}{N} \frac{f(\hat{X}_T^S | \hat{X}_{T-1}^i)}{f(\hat{X}_T^S | \hat{Y}_{1:T-1})} \delta(x_{T-1} - \hat{X}_{T-1}^i) \\
 &\approx \sum_{i=1}^N \frac{f(\hat{X}_T^S | \hat{X}_{T-1}^i)}{\sum_{j=1}^N f(\hat{X}_T^S | \hat{X}_{T-1}^j)} \delta(x_{T-1} - \hat{X}_{T-1}^i)
 \end{aligned}$$

So \hat{X}_{T-1}^S has probability of $\frac{f(\hat{X}_T^S | \hat{X}_{T-1}^i)}{\sum_{j=1}^N f(\hat{X}_T^S | \hat{X}_{T-1}^j)}$ to be equal to \hat{X}_{T-1}^i . The same

concept applies when sampling from \hat{X}_{T-2}^S , and so on.

Algorithm: particle smoother

1. Do particle filter first to get the filtering samples $\{\hat{X}_k^i : i=1, \dots, N, k=1, \dots, T\}$.
2. Resample $\hat{X}_T^S = \hat{X}_T^j$ with probability $\frac{1}{N}$
3. Compute $w_{T-1}^j = f(\hat{X}_T^S | \hat{X}_{T-1}^j)$ for $j=1, \dots, N$
4. Resample $\hat{X}_{T-1}^S = \hat{X}_{T-1}^j$ with probability $\frac{w_{T-1}^j}{\sum_{i=1}^N w_{T-1}^i}$
5. Repeat Step 3-4 for $k=T-1, T-2, \dots, 0$ to get $\{\hat{X}_T^S, \hat{X}_{T-1}^S, \dots, \hat{X}_0^S\} \sim f(x_{1:T} | \hat{Y}_{1:T})$
6. Repeat Step 2-5 M times to get M time-history samples of $f(x_{1:T} | \hat{Y}_{1:T})$.

Remarks:

1. Since the particle smoother is based on the particle filter results, so it is also asymptotically correct and not good for high dimensional X_k .
2. M is better not to be large than N since otherwise it seems like getting something from nothing.
3. After we obtain samples $\{\hat{X}_{1:T}^{S,i} : i=1, \dots, M\}$, $E[g(X_{1:T}) | \hat{Y}_{1:T}]$ can be estimated

as

$$E[g(X_{1:T}) | \hat{Y}_{1:T}] \approx \frac{1}{M} \sum_{i=1}^M g(\hat{X}_{1:T}^{S,i})$$