

Particle Filters

Alex Pan and Alec Kosik

Abstract

Particle Filters allow us to model non-linear and non-Gaussian distributions through Monte Carlo methods. We first discuss HMM, and then discuss a recursive formulation of the distribution of $p(x_t|y_{1:n})$. We then discuss how this leads well to an algorithmic approach. Finally, we investigate Kalman Filters, Importance Sampling, and SIR. We discuss the motivations for and relevant problems with each approach.

1 Hidden Markov Models

A **Hidden Markov Model** (HMM) is defined as a Markov process with states that cannot be observed, but gives outputs that are dependent on the hidden state of the system. If we let the random variable X_i represent the i -th state of the system, and the random variable Y_i represent the i -th output of the system, we have:

$$X_1 \sim \mu(x_1) \text{ and } X_n|(X_{n-1} = x_{n-1}) \sim f(x_n|x_{n-1}) \quad (1)$$

$$Y_n|(X_n = x_n) \sim g(y_n|x_n) \quad (2)$$

where X_1 is defined by some initial conditions $\mu(x_1)$. Note how f obeys the Markov property, so the current state x_n is dependent only on the previous state x_{n-1} . Also note that g shows that the current observation y_n is only dependent on the current state x_n .

Let $x_{1:n}$ represent the ordered n-tuple of hidden states (x_1, \dots, x_n) where x_i is the i -th hidden state. Similarly, let $y_{1:n}$ represent the ordered n-tuple of observations (y_1, \dots, y_n) . There are a few important distributions that follow directly from our HMM setting.

$$\begin{aligned} p(x_{1:n}) &= p(x_1)p(x_2|x_1)p(x_3|x_2, x_1) \dots p(x_n|x_{n-1}, \dots, x_1) \\ &= p(x_1)p(x_2|x_1)p(x_3|x_2) \dots p(x_n|x_{n-1}) \\ &= \mu(x_1) \prod_{i=2}^n f(x_i|x_{i-1}) \end{aligned} \quad (3)$$

$$\begin{aligned} p(y_{1:n}|x_{1:n}) &= p(y_1|x_{1:n})p(y_2|x_{1:n}, y_1)(y_3|x_{1:n}, y_2, y_1) \dots (y_n|x_{1:n}, y_{n-1}, \dots, y_1) \\ &= p(y_1|x_1)p(y_2|x_2) \dots (y_n|x_n) \\ &= \prod_{i=1}^n g(y_i|x_i) \end{aligned} \quad (4)$$

2 Particle Filters

2.1 General Goals

There are a few different uses of the particle filtering method. We could model the distribution $p(x_{1:n}|y_{1:n})$, the distribution $p(x_n|y_{1:n})$, or even the expected value of functions defined on these distributions. Here, our main goal is to use our observations to determine the hidden state of the system. That is, we want to figure out $p(x_n|y_{1:n})$. Recognizing that this is a marginal is useful because then we can just say:

$$p(x_n|y_{1:n}) = \int p(x_{1:n}|y_{1:n}) dx_{1:n-1} \quad (5)$$

Setting our goal as a marginal of $p(x_{1:n}|y_{1:n})$ is helpful because we can refactor the full conditional distribution to see how it comes from our prior distributions derived in (3) and (4). Refactoring and using our initial HMM conditions, we can show:

$$\begin{aligned}
p(x_{1:n}|y_{1:n}) &= \frac{p(x_{1:n}, y_{1:n})}{p(y_{1:n})} \\
&= \frac{p(x_{1:n})p(y_{1:n}|x_{1:n})}{p(y_{1:n})} \\
&= \frac{p(x_{1:n})p(y_{1:n}|x_{1:n})}{\int p(x_{1:n}, y_{1:n})dx_{1:n}} \\
&= \frac{\mu(x_1) \prod_{i=2}^n f(x_1|x_{i-1}) \prod_{i=1}^n g(y_i|x_i)}{\int \mu(x_1) \prod_{i=2}^n f(x_1|x_{i-1}) \prod_{i=1}^n g(y_i|x_i)dx_{1:n}}
\end{aligned} \tag{6}$$

Thus,

$$p(x_n|y_{1:n}) = \int \frac{\mu(x_1) \prod_{i=2}^n f(x_1|x_{i-1}) \prod_{i=1}^n g(y_i|x_i)}{\int \mu(x_1) \prod_{i=2}^n f(x_1|x_{i-1}) \prod_{i=1}^n g(y_i|x_i)dx_{1:n-1}} dx_{1:n-1} \tag{7}$$

This looks really messy, but it's helpful to see that this is just integration of functions we already know: f and g . In other words, we already have all the information we need to arrive at the answer.

2.2 A Recursive Formulation

While we now have a nice equation, it would be really inefficient if we had to calculate this whole equation every single iteration of our system. Instead, we can formulate this recursively and see how some of the values carry through. This recursiveness is described in two parts: the update equation and the prediction equation. We show the update equation first:

$$p(x_t|y_{0:t}) = \frac{p(x_t|y_{0:t-1})p(y_t|x_t)}{p(y_t|y_{0:t-1})} \tag{8}$$

The derivation will be broken into 2 parts. First, we will show that (8) holds iff $p(x_t, y_{0:t}) = p(x_t, y_{0:t-1})p(y_t|x_t)$ holds. Then we will show that $p(x_t, y_{0:t}) = p(x_t, y_{0:t-1})p(y_t|x_t)$. We start by applying Bayes' rule to the left side of (8) and then moving the denominator to the right side:

$$p(x_t, y_{0:t}) = \frac{p(x_t|y_{0:t})p(y_t|x_t)}{p(y_t|y_{0:t-1})}p(y_{0:t})$$

Applying Bayes' rule to $p(x_t|y_{0:t})$:

$$\begin{aligned}
&= \frac{p(x_t, y_{0:t})p(y_t|x_t)}{p(y_t|y_{0:t-1})p(y_{0:t-1})}p(y_{0:t}) \\
&= \frac{p(x_t, y_{0:t})p(y_t|x_t)}{p(y_{0:t})}p(y_{0:t}) \\
&= p(x_t, y_{0:t})p(y_t|x_t)
\end{aligned}$$

Second, we show $p(x_t, y_{0:t}) = p(x_t, y_{0:t-1})p(y_t|x_t)$. Let's start by refactoring the right-side.

$$\begin{aligned}
p(x_t, y_{0:t-1})p(y_t|x_t) &= p(x_t|y_{0:t-1})p(y_{0:t-1})\frac{p(x_t, y_t)}{p(x_t)} \\
&= p(x_t|y_{0:t-1})p(y_{0:t-1})\frac{p(x_t|y_t)p(y_t)}{p(x_t)}
\end{aligned}$$

Using the fact that: $p(x|yz) = \frac{p(y)p(z)}{p(yz)} \frac{p(x|y)p(x|z)}{p(x)}$

Let: $x = x_t, y = y_{0:t-1}, z = y_t$,

We can now see that: $p(x_t|y_{0:t-1}, y_t) = \frac{p(y_{0:t-1})p(y_t)}{p(y_{0:t-1}, y_t)} \frac{p(x_t|y_{0:t-1})p(x_t|y_t)}{p(x_t)}$

Thus,

$$\begin{aligned} p(x_t|y_{0:t-1})p(y_{0:t-1}) \frac{p(x_t|y_t)p(y_t)}{p(x_t)} &= p(x_t|y_{0:t})p(y_{0:t-1}) \\ &= p(x_t, y_{0:t}) \end{aligned}$$

We now show how to derive the prediction equation:

$$p(x_t|y_{0:t-1}) = \int f(x_t|x_{t-1})p(x_{t-1}|y_{0:t-1})dx_{t-1} \quad (9)$$

If we approach this as a marginal by integrating out x_{t-1} we normally would formulate this as:

$$\begin{aligned} p(x_t|y_{0:t-1}) &= \int p(x_{t-1:t}|y_{0:t-1})dx_{t-1} \\ &= \int p(x_t|x_{t-1}, y_{0:t-1})p(x_{t-1}|y_{0:t-1})dx_{t-1} \end{aligned} \quad (10)$$

We can make use of the Markov property to see:

$$p(x_t|x_{t-1}, y_{0:t-1}) = p(x_t|x_{t-1}) = f(x_t|x_{t-1}) \quad (11)$$

———Some more words explaining what exactly just happened should go here.———

The issue now is that this is sometimes intractable. If functions are nice, like guassians or linear distributions, this is integrable and we can arrive at a solution analytically. However, this is not always the case, which is why we sometimes need to derive an answer from numeric approximations.

3 Kalman Filter

This needs to be re-done because I've moved the discussion of the recursive form to the earlier section. Former code has been commented out.

4 Sequential Monte Carlo Methods

So what if our distribution isn't nice? We will use Monte Carlo Methods to approximate Eq. (5) and (6):

$$\text{Prediction step: } \pi_{x_{1:t}|y_{1:t-1}}(x_{1:t}|y_{1:t-1}) = f(x_t|x_{t-1})\pi_{x_{1:t-1}|y_{1:t-1}}(x_{1:t-1}|y_{1:t-1}) \quad (12)$$

$$\text{Update step: } \pi_{x_{1:t}|y_{1:t}}(x_{1:t}|y_{1:t}) \propto g(y_t|x_t)\pi_{x_{1:t}|y_{1:t-1}}(x_{1:t}|y_{1:t-1}) \quad (13)$$

5 SIS

1. Initialization:

For $i = 1, \dots, N$:

Sample $x_0^{(i)} \sim \mu(x_0)$

Assign weights $\tilde{w}_0^{(i)} = g(y_0|x_0^{(i)})$

Normalize weights $w_0^{(i)} = \frac{\tilde{w}_0^{(i)}}{\sum_{i=1}^n \tilde{w}_0^{(i)}}$

2. Importance Sampling:

For $t = 1, \dots, T$:

Sample $x_t^{(i)} \sim f(x_t|x_{t-1}^i)$

Assign weights $\tilde{w}_t^{(i)} = g(y_t|x_t^{(i)})$

Normalize weights $w_t^{(i)} = \frac{\tilde{w}_t^{(i)}}{\sum_{i=1}^n \tilde{w}_t^{(i)}}$

3. Return $\{x_T^{(i)}, w_T^{(i)}\}_{i=1}^N$

To algorithmically find our target distribution, we sample from something simple like the Uniform or Gaussian distribution. We then weight these samples using our observations and g . Essentially, we iterate between (7) and (8), sampling new points and updating our weights each time. While this sounds okay, it turns out that we never use SIS in real life because of something called the “degeneracy problem”.

6 Degeneracy Problem

When we are sampling from our reweighted distribution, it is possible that certain points are assigned very low weight. In future iterations, this is likely to compound, making their weights essentially negligible. This manifests as a handful of points having nearly all the weight while the rest of our points have “degenerated”. This is a problem because our distribution essentially is now defined by just a few of the original points we started with, which is not enough to get an accurate representation of our target distribution.

7 Bootstrap (SIR)

1. Initialization:

For $i = 1, \dots, N$:

Sample $x_0^{(i)} \sim \mu(x_0)$

Assign weights $\tilde{w}_0^{(i)} = g(y_0|x_0^{(i)})$

Normalize weights $w_0^{(i)} = \frac{\tilde{w}_0^{(i)}}{\sum_{i=1}^n \tilde{w}_0^{(i)}}$

2. Importance Sampling:

For $t = 1, \dots, T$:

Sample $x_t^{(i)} \sim f(x_t|\tilde{x}_{t-1}^i)$

Assign weights $\tilde{w}_t^{(i)} = g(y_t|x_t^{(i)})$

Normalize weights $w_t^{(i)} = \frac{\tilde{w}_t^{(i)}}{\sum_{i=1}^n \tilde{w}_t^{(i)}}$

Resample $\tilde{x}_t^{(i)}$ from $x_t^{(i)}$ according to the weight distribution with replacement.

3. Return $\{x_T^{(i)}, w_T^{(i)}\}_{i=1}^N$

While this addresses the degeneracy problem, it turns out that resampling has its own issue which we call the “sampling impoverishment problem”. This comes about because our higher weighted points are more likely to be drawn multiple times. Instead of having many low-weight points spread out across the state-space—in the case of SIS and the degeneracy problem—we will have lots of evenly-weighted points concentrated around areas where high-weight points originally occurred. Intuitively, the degeneracy problem leaves you with a diverse scattering of point masses with a few ill-defined peaks while the sampling impoverishment problem leaves you with *only* well defined peaks. There are methods that address both of these issues but they require fancy sounding things like the “Epanechnikov Kernel”, which is outside the scope of our discussion. It turns out that SIR is generally good enough and is the algorithm that people generally use.