

More Statistics for IIB Project

Name: Zhenhua Lin
Crsid: zhl24
College: Fitzwilliam College

February 22, 2024

Abstract

1 Understanding of the Big Story

In Michaelmas term (work documented in the previous report), what we did was building a framework for generating and inferring state space models driven by Levy process, NVM process as the most discussed example. In the inference part, we assumed that the exact driving process is known with full details, including the jump sizes and times, making the rest of the system Gaussian (hence called a conditionally Gaussian system), and we implemented a Kalman filter to infer in such case. In this work, we first narrow down the assumption of knowing the exact process driving the system with full details to knowing the exact generator for the driving process, and use it as the proposal for a bootstrap particle filter which can be used alone or combined with a Kalman filter to form the marginalized particle filter to infer in the system. After that, we probably would further narrow down the assumption to knowing the kind of process and infer the parameters? And eventually no prior knowledge?

2 Particle Filter

Inference in stochastic system only has closed form solution for the posterior in some very limited cases, including the simple linear Gaussian model (system with linear dynamic equations i.e. not quadratic or sinusoidal and driven by Gaussian noise) and the general discrete hidden Markov model. In the vast majority of cases, non-linearity or non-Gaussianity render an analytic solution intractable.

For the Gaussian systems with non-linear dynamics, one classic solution is the extended Kalman filter. The principle is simple, by using Taylor expansion to have a linear approximation of the non-linear dynamics that is compatible with Kalman filter. This method, obviously, fails if the system has substantial non-linearity, or the system is driven by highly non-Gaussian noise. There are some other approaches based on approximation of the non-Gaussian distribution, such as via mixture of Gaussians or second order Taylor series, but all of them have limitations to the form of the posterior distributions.

We hence want a method that can achieve non-linear filtering with no requirements in the form of posterior distribution. And methods based on Monte Carlo satisfy such needs, but with a significant cost in the computational power. However, it is still quite a popular choice because of its generality that allows inference of full posterior distributions in general state-space models.

2.1 General Bayesian Filtering Problem

The main reference for this session is the paper about particle filter tutorial.

Here we study the most important sampling approach for hidden states, particle filtering. Particle filter does similar functions as Kalman filter, but the difference is that the particle filter can be applied to both linear

and non-linear dynamics. The filter enables representing state estimates with arbitrarily shaped probability distributions. The key questions lie in how to set up a particle filter for the inference problem, what would be the outputs and how would the outputs be fed to the system to provide effective inference.

As mentioned before, particle filter solves the problem of similar structure but compatible with non-linear dynamics:

1. Process model for state changes over time given noise and optional inputs:

$$\mathbf{x}_k = f_k(\mathbf{x}_{k-1}, \mathbf{u}_k, \mathbf{v}_{k-1}) \quad (1)$$

2. A measurement model of the form:

$$\mathbf{z}_k = \mathbf{h}_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{n}_k) \quad (2)$$

The main difference is that both \mathbf{f} and \mathbf{h} can be arbitrarily shaped distributions.

The particle filter is a Bayesian filter i.e. estimation is performed using Bayesian theory, allowing for estimating a state by combining a statistical model for measurement(likelihood) with a prior probability. The estimation is solved in a similar form to Kalman filter that applies recursive predict-update cycle.

The main steps for such Bayesian filter have the same formulae (both Kalman and Particle filters belong to this class):

1. The **predict step** basically combines the posterior estimate of the previous state with the process model to estimate the current state without observation, via a marginalizing integral. The specific formula can be obtained via marginalizing the posterior estimate of the current state given previous observations. This essentially makes use of the hidden Markov structure.

$$p(\mathbf{x}_k | \mathbf{z}_{1:k-1}) = \int p(\mathbf{x}_k | \mathbf{x}_{k-1}) p(\mathbf{x}_{k-1} | \mathbf{z}_{1:k-1}) d\mathbf{x}_{k-1} \quad (3)$$

This formula can be understood as the forward simulation from the previous state posterior prediction making use of the system hidden state transition dynamics. The weighted sum of the forward simulations varied in each initial hidden state is the final prediction of the next hidden state, i.e. Bayesian prediction.

2. The **update step** is basically to estimate the current state after a current observation by combining the prediction posterior and emission mechanism (simply \mathbf{h}) and also a marginalisation constant via the Bayes theorem. Essentially negative feedback of the prediction using the observation at the corresponding time position.

$$p(\mathbf{x}_k | \mathbf{z}_{1:k}) = \frac{p(\mathbf{z}_k | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{z}_{1:k-1})}{p(\mathbf{z}_k | \mathbf{z}_{1:k-1})} \quad (4)$$

Where $p(\mathbf{z}_k | \mathbf{x}_k) = \mathbf{h}_k(\mathbf{x}_k, \mathbf{u}_k, \mathbf{n}_k)$ and $p(\mathbf{x}_k | \mathbf{z}_{1:k-1})$ is just from the last step.

The more difficult one is the denominator, the normalizing term, which is also updated in every update step by combining the predict posterior probability with the emission mechanism. This is important in some cases with dynamic structure e.g. noise parameters as the hidden states to be predicted.

$$p(\mathbf{z}_k | \mathbf{z}_{1:k-1}) = \int p(\mathbf{z}_k, \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{z}_{1:k-1}) d\mathbf{x}_k \quad (5)$$

Pay attention to the time index in the formulae which provide very important intuition but may be messed up.

The integrals in the predict and update steps for Bayesian filters as shown in equation (3) and (5) are only analytically solvable under strong assumptions : finite dimensional discrete state variables or linear models with Gaussian pdfs. What a particle filter does is basically approximating the posterior pdf by a discrete pdf such that there are minimal restrictions in the models involved. The optimal Bayesian solution is

approximated by a sum of weighted dirac delta samples. Particle filtering is essentially an approximation of the whole probability space via a combination of multiple pointwise evaluation carried by the offsetted delta function corresponding to the true place.

$$p(\mathbf{x}_{0:k}|\mathbf{z}_{1:k}) \sim \sum_{i=1}^{N_s} w_k^i \delta(x_{0:k} - x_{0:k}^i) \quad (6)$$

The intuition behind particle filtering is that, we have multiple paths or particles each of which represents a specific evaluation in the whole probability space. Therefore, if we can have an infinite number of particles, we could ideally have a complete evaluation of the whole probability space. (This summation formula is not especially meaningful just to point out that this is an approximation for the true space).

Where $\{w_k^i, \mathbf{x}_{0:k}^i\}_{i=1}^{N_s}$ is the set containing N_s sample sequences and weights, with each weight representing the relative importance of each of sample sequences, and the sum of weights is 1. Samples associated with high weights are believed to be closer to the true state sequence than samples associated with low weights. The integral of the sum of delta functions i.e. equation (6) is hence the summation of the weights and equal to one, thus a valid probability distribution approximation.

The discrete approximation of a continuous pdf turns intractable integrals into summations over N_s samples, which are often referred to as the particles.

2.2 Algorithms for Particle Filtering (The following is mainly just for the algorithms, specific derivation should check professor Godsill's overview paper)

The discrete approximation of a continuous pdf turns intractable integrals into summations over N_s samples. The samples are usually referred to as particles hence the name particle filter. The advantages of representing the posterior by a set of weighted particles include (i) the ability to represent arbitrarily shaped pdfs (assuming enough samples) and (ii) minimal restrictions on the process and measurement models. This combination of advantages is one of the main reasons for the popularity of the particle filter. The most obvious drawbacks are (i) the lack of expressiveness in case the number of particles is too low and (ii) the increased computational costs compared to the usual inference methods.

The sample-based approximation comes with an obvious challenge. The posterior pdf that must be estimated is unknown hence sampling from it is impossible. Samples must therefore be drawn from another distribution instead. This distribution, referred to as importance density (sometimes called proposal density) will be denoted q . The weights compensate for the fact that samples are drawn from the importance density q rather than the posterior pdf. Any function that is positive where the posterior is positive can be used as importance density (sequential importance sampling).

Particles are basically just samples with weights in the filtering scheme.

2.2.1 The Most Basic Particle Filter: Sequential Importance Sampling (SIS)

The simplest particle filter algorithm is hence just sequential importance sampling, and we assume q to be a generator function we have and p to be the desired distribution. In the filtering scheme, both q and p refer to the transition density function in the system.

1. We first define the number of particles we want for the algorithm: N_s , and we draw N_s samples or particles from the prior distribution $q(\mathbf{x}_0)$ and assign the initial weights for each sample as :

$$w_0^i \propto \frac{p(\mathbf{x}_0^i)}{q(\mathbf{x}_0^i)} \quad (7)$$

The weights are always defined as the ratio of the posterior hidden states probability given observations (none for prior) to the importance probability of observing the data sequence from the generator probability.

2. Then in each step, we propagate the particles in last step: $\{x_{k-1}^i, w_{k-1}^i\}_{i=1}^{N_s}$ as the current observations \mathbf{z}_k arrive, to obtain the new particles $\{x_k^i, w_k^i\}_{i=1}^{N_s}$.

We first draw samples from the proposal distribution q , conditional on the last trajectory of the particle and the current observation (need Bayesian expansion to evaluate, note that this is not the transition density):

$$x_k^i \sim q(x_k^i | x_{k-1}^i, \mathbf{z}_k) \quad (8)$$

Next we assign weights to the samples to make them particles, which equal to the ratio of the exact posterior distribution to the proposed one, but note that these are evaluated over the whole trajectory for each particle, not just the instantaneous values:

$$w_k^i \propto \frac{p(\mathbf{x}_{0:k}^i | \mathbf{z}_{1:k})}{q(\mathbf{x}_{0:k}^i | \mathbf{z}_{1:k})} \quad (9)$$

But note that this equation is not trivial to evaluate directly, and instead, we solve it by using a similar recursive update form by rewriting the equation as:

$$w_k^i \propto w_{k-1}^i \frac{p(z_k | x_k^i) p(x_k^i | x_{k-1}^i)}{q(x_k^i | x_{k-1}^i, z_k)} \quad (10)$$

with the numerator terms being the emission and transition density functions, and the denominator term being the proposal density function that needs Bayesian expansion to evaluate (samples generated not just based on the last states but also the current observation, two inference paths). Use simple sum-to-one constraint to have the exact expression from the variation part.

Then we have the particles for the current time step: $\{x_k^i, w_k^i\}_{i=1}^{N_s}$

3. The particles give us a probability distribution of the hidden states at the time step according to formula (6) i.e.

$$p(\mathbf{x}_{0:k} | \mathbf{z}_{1:k}) \sim \sum_{i=1}^{N_s} w_k^i \delta(x_{0:k} - x_{0:k}^i) \quad (11)$$

with the sum of weight history of a particle equal to the probability of the hidden states being the trajectory of a specific particle (discrete distribution approximation over some particular sequences). Then we could do something like ML estimate to infer the hidden states of the noisy signal using non-linear dynamics.

An additional thing to note is that we often use the normalized weights for the particles instead. So the weight of the i th particle now becomes:

$$\hat{w}^{(i)} = \frac{w^{(i)}}{\sum_j w^{(j)}} \quad (12)$$

This has a simple derivation:

1. Say the target of estimate is some random function $h(x)$, such that we want to know the expectation:

$$E(h(x)) = \int h(x) p(x) dx \quad (13)$$

2. Again by importance sampling, we have the estimation formula:

$$E(h(x)) = \int h(x) p(x) \frac{q(x)}{q(x)} dx \sim \frac{1}{N} \sum_i \frac{p(x^{(i)})}{q(x^{(i)})} h(x^{(i)}) = \frac{1}{N} \sum w^{(i)} h(x^{(i)}) \quad (14)$$

3. To use the normalized weights, we need an additional relation:

$$1 = \int p(x)dx = \int p(x) \frac{q(x)}{q(x)} dx \sim \frac{1}{N} \sum_i \frac{p(x^{(i)})}{q(x^{(i)})} = \frac{1}{N} \sum_i w^{(i)} \quad (15)$$

The key to note is that we have the similarity:

$$N \sim \sum_i w^{(i)} \quad (16)$$

4. Finally, we put the expression for N in (16) into (14) to obtain the estimation formula in the normalized weights:

$$E(h(x)) \sim \sum_i \frac{w^{(i)}}{N} h(x^{(i)}) = \sum_i \frac{w^{(i)}}{\sum_j w^{(j)}} h(x^{(i)}) = \sum_i \hat{w}^{(i)} h(x^{(i)}) \quad (17)$$

Note that the important thing is that h can be anything we want to estimate. The most straight forward example is the mean $E(x)$, which would be just the weighted sum of the particles, way more reasonable than the simple average.

2.2.2 Bootstrap Particle Filter: SIS with Resampling (For specific derivation, check the PF paper by Professor Godsill, the following is just the algorithm)

However, the simplest SIS algorithm has several problems and modifications are needed to make it an effective particle filter algorithm. The most famous problem is the weight degeneracy problem. Weight degeneracy often occurs in scenarios where the likelihood function sharply contrasts the prior or proposal distribution. In such cases, particles that are not close to the high probability regions of the likelihood function end up with negligible weights. Another common cause is the accumulation of errors over time in the estimated state, which leads to a situation where most particles fail to predict the observed data accurately, resulting in their weights becoming insignificant. Note that this doesn't have much relation to the MH problem of acceptance probability tending to 0.

The solution to that is by implementing a sequential importance sampling algorithm with resampling for the particles.

The problem is defined as below:

After a few iterations one particle weight will be very close to one and all other particle weights will be almost zero. The consequences of the degeneracy problem are (i) almost all computational effort will be put into computations related to particles that have negligible or no contribution to the overall estimate and (ii) the number of effective particles is only one. The latter greatly limits the performance and expressiveness of the filter. The expressiveness is limited since a single particle can only represent one point in the state space rather than pdfs of arbitrary shapes. Performance is poor since the particle filter will diverge. Divergence occurs in case state estimation errors increase over time and are unacceptably large. Once diverged, the filter fails to "follow" the true state or reduce estimation errors to acceptable values.

The solution to this problem is resampling. In the most basic resampling, say we have decided to have N trajectories, using the normalized weights as the probabilities of selection of the particles, we draw N particles with replacement. Therefore, trajectories with small importance weights are eliminated, whereas those with large importance weights are replicated. The normalized importance weights are set to $1/N$ after the resampling step.

Therefore we have an algorithm for the basic particle filter with resampling:

1. To start the algorithm, we need several inputs, including the number of trajectories we want N, initial particles $\{X_0\}$, observation data, emission function from samples to observation g, and the transition function for particles q (we denote f as the actual transition function). In bootstrap particle filter, $f=q$ directly.

Note also that we need to compute the initial weights for the initial set of particles as in (18) which shows the weight for the ith particle without normalization yet. Note that we are only interested in the

contribution from the i th particle, so the condition is on the i th hidden state but not the whole hidden states in the beginning, where π_0 is the initial state distribution. We need the initial probabilities for the weights. Note that we need to normalize the weights after the computation. **Note what weights except the normalised weights should all be in the log domain! Normalisation should also be via logsumexp.**

$$w_0^{(i)} = \frac{g(y_0|x_0^{(i)})\pi_0(x_0^{(i)})}{q_0(x_0^{(i)}|y_0)} \quad (18)$$

The transition function should take the particles and return the propagated particles to the next time step through the system. The emission function should take the observation and particles at the same time step to return the corresponding emission probability.

The initial particles should be sampled from some initial prior.

2. Then we run the algorithm through each time step in the observation data to infer the hidden states. Note that SIS method applies to time series data.
3. In each iteration after initialization, we perform **resampling** first for the particles from the last step, based on the normalized weights $\{\hat{w}_{t-1}\}$ which are taken directly as the selection probabilities. We draw N particles with replacement according to the normalized weights. After resampling, we set the normalized weights of the resampled particles evenly to $\frac{1}{N}$.

Then we perform the **propagation** step for the resampled particles based on the transition kernel function provided (return new proposed samples directly). For the i th trajectory (often no dependence on observation y , just the previous particles):

$$x_t^{(i)} \sim q(x_t^{(i)}, y_t) \quad (19)$$

Finally we **update and normalize the weights**. The unnormalized weights are computed according to equation (20). Note that we need the kernel to compute the transition probability in the general case. Remember to normalize the weights. w_{t-1} is basically the uniform weights passed from the last step.

$$w_t^{(i)} = w_{t-1}^{(i)} \frac{g(y_t|x_t^{(i)})f(x_t^{(i)}|x_{t-1}^{(i)})}{q_t(x_t^{(i)}|x_{t-1}^{(i)}, y_t)} \quad (20)$$

But for the **weight update in the bootstrap particle filter**, $f=g$ and we need only the emission probability as shown in (21). Remember also to normalize the weights.

$$w_t^{(i)} = w_{t-1}^{(i)} g(y_t|x_t^{(i)}) \quad (21)$$

2.3 Particle Filter for Levy State Space Model

We first implement a basic bootstrap particle filter to the Levy state space model. Hence, what we need are the transition function for propagating the particles, and the emission function for computing the probability of a particle given an observation.

First, revise the formula for the Levy state space model driven by an NVM process, the Normal Gamma process to be more specific, as shown in equation (22).

$$\mathbf{X}(t) = \mathbf{e}^{\mathbf{A}(t-s)} \mathbf{X}(s) + \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} [\mu_w Z_i + \sigma_w \sqrt{Z_i} U_i] \mathbf{e}^{\mathbf{A}(t-V_i)} \mathbf{h} \quad (22)$$

We further evaluate it into the form that is convenient for simulation, with Δt being the simulation time step size, $\{Z_i, V_i\}$ being the jump sizes and times for the subordinator process, and σ_w, μ_w being the NVM parameters. Note that such system has a fixed transition function given uniform simulation time step, but the system exponent in the noise has to be calculated for each jump and evaluation point. We compute the

new sample by simulate the driving process first and then finidng the corresponding jumps in the interval using formula (22). **Note that in our simulation scheme, we divide the time axis into many dt intervals. So the jump times we obtained are essentially $\tau_i = V_i - t$, thus whats inside the bracket is just minus of the offseted jump time in the dt interval $-\tau_i$. Hence, we dont need to bother with the exact time point of the system t.**

$$X_{n+1} = e^{\mathbf{A}\Delta t} X_n + \sum_{\frac{t-\Delta t}{T} \leq V_i \leq \frac{t}{T}} [\mu_W Z_i + \sigma_W \sqrt{Z_i} U_i] \mathbf{e}^{\mathbf{A}(t-V_i)} \mathbf{h} \quad (23)$$

The observations are again the true hidden states covered by Gaussian noise of some variance.

$$\mathbf{Y}_n = \mathbf{X}_n + \sigma \mathbf{U}_n \quad (24)$$

Then, what we need to do is to draw samples from the transition function $f(\mathbf{X}_{n+1}|\mathbf{X}_n)$ for forward simulation or sampling and define the observation function $g(\mathbf{Y}_n|\mathbf{X}_n)$ i.e. the weight function to run a bootstrap particle filter. Note that the transition function here probably does not have a closed form, and there is no need to define it, because we just need to simulate forwards to have new samples proposed. That is for the transition function, all we need is the exact driving process generator to simulate forwards by dt to obtain the jump sizes and times.

We can therefore understand particle filtering scheme as simulating multiple paths to try to reconstruct the real one. It is just having multiple simulations from zero, and using weights as the measure for matching of a simulation to the real one (constrained by weight which is higher for more likely simulations given the observations). It is something like the weighted sum of all simulated paths. The observation probability is shown in equation (25). Note that since additional normalisation step would be carried out, and the normalisation constant is the same for all samples since it depends on the data only, and as mentioned in the paper, that would be constant for fixed data filtering scheme, there is no need to include it in the function. We just need a score for matching.

$$g(\mathbf{Y}_n|\mathbf{X}_n) = N(\mathbf{X}_n, \sigma^2 \mathbf{I}) \sim \exp((\mathbf{Y}_n - \mathbf{X}_n)^T (\mathbf{Y}_n - \mathbf{X}_n) / 2\sigma^2) \quad (25)$$

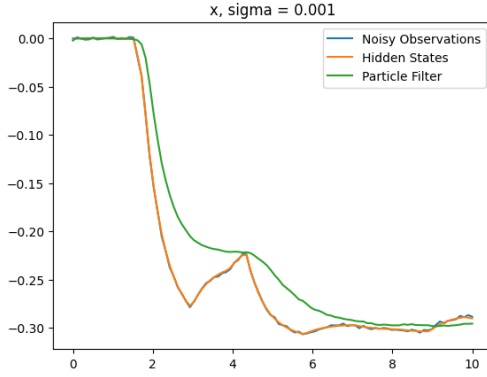
The main problem in building the particle filter transition function is that since it is just some random process generator, but we do not necessarily know the type of process or the exact parameters for such generator. For this case, we would just assume that we know exactly the type of process and parameters, by using the exact generator that generates the data, and leave such problem for latter discussion.

Note that in coding, always remember that sample generation always has the chance of returning empty array which could cause dimension mismatch.

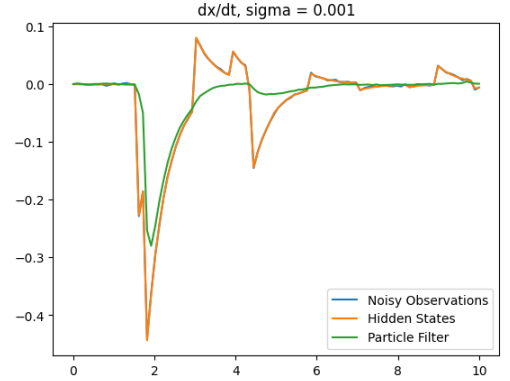
Hence, we have a simple implementation of particle filtering in the Levy state space model. Some example results are shown in figure 1 below.

3 Marginal Particle Filter for NVM State Space System

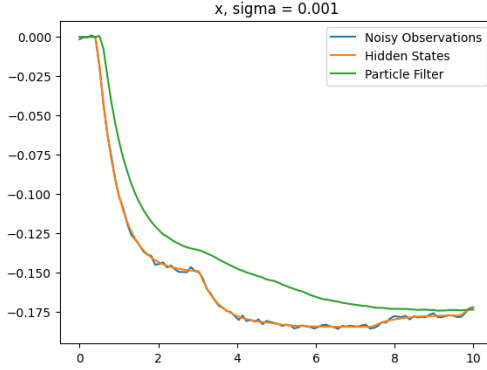
Marginal particle filter can be applied to problems where the linear and non-linear parts have clean separation. This is exactly the case of our system driven by NVM process, where if we condition on the non-linear states, the rest of the states are just simple Gaussian process. In marginal particle filter, instead of using particle filter to propose the complete dynamics of the system, we use it to propose only the non-linear part of states in each time step. And condition on the non-linear proposal, the rest of the states are just simple Gaussian which can be estimated by linear approach i.e. Kalman filter. Such clean separation is aligned with the nature of a system driven by the NVM process, which sequentially generates the non-linear subordinator process first and then generates the linear process conditional on the non-linear parts. The main technical difficulty lies in how to compute the weights of the particles in this case, and would be discussed in the following session with comparison with the previous pure particle filter case.



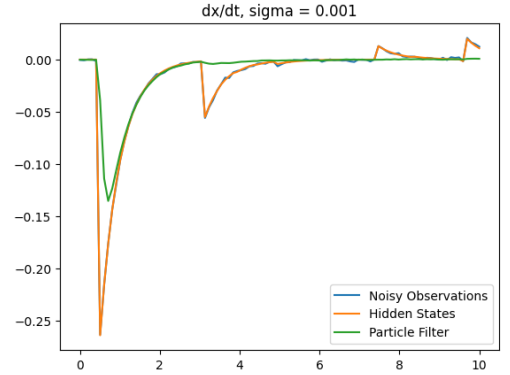
(a) Particle Filtering with 1000 Particles (x)



(b) Particle Filtering with 1000 Particles ($\frac{dx}{dt}$)



(c) Particle Filtering with 10000 Particles (x)



(d) Particle Filtering with 1000 Particles ($\frac{dx}{dt}$)

Figure 1: Particle Filtering Implementation Example Results

3.1 Derivation of Marginal Particle Filter Algorithm for Levy State Space System

We again begin by listing out the Levy state space dynamic model as shown in equation (26). Note that in the previous pure particle filtering case, the proposal generates the jump times and the NVM jump sizes (whats inside the bracket). The special thing is to note that that would include both the linear dynamics U_i and the non-linear dynamics Z_i . A seemingly appealing structure would be proposing the same NVM jump sizes again, and hence the weight function is $p(Y|\text{proposal})$ would be equivalent to $p(Y|X)$, the emission probability again, since the NVM jump sizes and times form deterministic transformation to the hidden state, and then conditional on NVM jump sizes we know the subordinator jumps also (information subset) and can have the mean and covariance of the conditional system for us to run a Kalman filter to give the final output. But that would not be a valid marginal particle filter, and it is in fact doing more or less the same thing as the standard particle filter statistically, but only with an additional Kalman filtering step which introduces more uncertainty, because that proposal is still proposing the combined states.

$$X_{n+1} = e^{\mathbf{A}\Delta t} X_n + \sum_{\frac{t-\Delta t}{T} \leq V_i \leq \frac{t}{T}} [\mu_W Z_i + \sigma_W \sqrt{Z_i} U_i] \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} \quad (26)$$

Therefore, what a valid proposal should look like should be proposing only the subordinator jump sizes Z_i and times, leaving the Gaussian part in the NVM jump U_i to be inferred by a Kalman filter, such that we can make use of the nature of Kalman filter as the optimal inference solution to improve the inference results. **(The key thing to realize is that, the nature of the SDE system depends completely**

on the underlying driving process. The separation between linear and non-linear parts of the system lies in the separation property in the NVM process, but not the system.)

$$\mathbf{m} = \mu_W \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} Z_i \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} \quad (27)$$

$$\mathbf{S} = \sigma_W^2 \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} (\sqrt{Z_i} \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h}) (\sqrt{Z_i} \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h})^T \quad (28)$$

Then, we can define the weights for particles using $p(\text{observation}|\text{proposal}) = p(Y|Z_i, \tau_i)$. Note that these conditions would be equivalent to the mean and variance in the conditionally Gaussian term, since the transformation between them is deterministic as shown in (27) and (28). The weight would hence be calculated by $p(Y|\mathbf{m}, \mathbf{S})$ and that is in fact the marginal likelihood in the conditional Gaussian system, and we can estimate it via Kalman filter directly. The conditionally Gaussian system structure we would be working on is in the form of standard hidden Markov state transition and emission structure, with $V_n \sim N(\mathbf{m}, \mathbf{S})$. Note that the transition noise is a biased Gaussian noise here, simply add the bias to the deterministic transition and the noise would become zero mean, and the system would be compatible with the classical zero mean noise Kalman structure again.

$$X_{n+1} = F X_n + V_n \quad (29)$$

$$Y_n = H X_n + W_n \quad (30)$$

The marginal weights according to the Kalman filter notes by professor Godsill should be calculated as shown in equation (31).

$$p(y_t|y_{1:t-1}) = N(\hat{X}_{t|t-1}, S_{x_t} + \sigma_n I) \quad (31)$$

This is in fact simply taking expectation wrt all terms using the simple hidden state emission linear equation. The particles here are just the conditional Gaussian parameters computed using the gamma jump sizes and times and the system structure. Inference of the hidden states is done by running a Kalman on the conditionally Gaussian system. σ_n is just the observation noise covariance matrix.

Therefore, a complete marginalised particle filtering algorithm would look like:

1. We first define the number of particles and the proposed system for generating the subordinator processes for particle filtering. Give also some initialisation to particle filter and Kalman filter.
2. In each iteration, we first use the proposed system to generate the non-linear states, subordinator jump sizes and times, as the particles. These could be transformed deterministically and directly into the Gaussian mean and variance in the conditionally Gaussian system, like in equation (27) and (28). Note that this generation is independent in different iterations. Also we don't have enough information to compute the weights just based on the non-linear states here.
3. Each particle would represent a conditionally Gaussian system (deterministic transformation between subordinator process and conditional Gaussian parameters). We then run a Kalman filter for each particle for one transition and correction step to infer the hidden states. The marginal probability can be computed directly in the Kalman correction step, which is exactly the weight for the corresponding particle. We then normalize the weights (always in the log domain!) for resampling.
4. Then, we can eventually do resampling of particles. Note that, not only just the particles need to be resampled, but also the inferred states by Kalman filter. We need to re-run the Kalman filter for each resampled particle to produce the new inference results (2 runs!). The new weights should also be calculated.
5. Finally, we can combined the inference results using the new samples. That should be done before resampling using the true weights and brand-new predictions, since that would correspond to the lowest Monte Carlo error, though after resampling using the new weights computed (**not the uniform weights re-assigned! Bad performances!**) still has quite good performances.

3.2 Extended State in Marginalised Particle Filter for Parameter Estimation

This is a method mentioned in the Levy state space model paper, that is particularly for the marginalised particle filtering method, based on the bias term of the Gaussian noise in the conditionally Gaussian system, which carries information about μ_w in an NVM process so we are essentially inferring the parameter while we are inferring the mean of the conditional Gaussian distribution.

We again start from the 3 important equations for derivation of such method in the NVM system. What this mainly implies to the non-linear state proposal is that, we now no longer have the full information to deterministically transform the jump sizes and times to the conditional Gaussian distribution mean parameter (equation (44) can no longer be used here), and the variance can still be deduced in the same way as before, since now μ_w is taken away as a state.

This changes the hidden Markov structure, especially the key driving conditional Gaussian distribution, since the mean is now taken out to be used as an external state. The driving Gaussian is now a zero mean one with the same covariance structure as before.

$$X_{n+1} = e^{\mathbf{A}\Delta t} X_n + \sum_{\frac{t-\Delta t}{T} \leq V_i \leq \frac{t}{T}} [\mu_w Z_i + \sigma_w \sqrt{Z_i} U_i] \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} \quad (32)$$

$$\mathbf{m} = \mu_w \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} Z_i \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} \quad (33)$$

$$\mathbf{S} = \sigma_w^2 \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} (\sqrt{Z_i} \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h}) (\sqrt{Z_i} \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h})^T \quad (34)$$

We can see how to include the extended state immediately once we formulate the conditionally Gaussian hidden Markov system with bias extracted to make the noise zero-mean, as shown in equation (35), with $V_n \sim N(0, \mathbf{S})$. The mean of the noise is extracted as the additional constant. This bias naturally carries the information about μ_w , and can hence be used for estimating μ_w in the NVM process.

$$X_{n+1} = e^{\mathbf{A}\Delta t} X_n + \mu_w \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} Z_i \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} + V_n \quad (35)$$

Since μ_w is no longer a part of the conditional Gaussian distribution, we need to define our new state space system with an extended state of μ_w . Hence the state vector is now shown in equation (46). Note that we need not go back to our exact differential equation to come up with a new \mathbf{A} , and what we need to do is just to define a new hidden Markov model with the new state vector.

$$\hat{\mathbf{X}} = \begin{bmatrix} \mathbf{X} \\ \mu_w \end{bmatrix} = \begin{bmatrix} x \\ \frac{dx}{dt} \\ \mu_w \end{bmatrix} \quad (36)$$

Here, we outline the marginalised particle filtering scheme with an augmented state.

1. First, we have our particle filter proposing the jump sizes and times of the subordinator process as usual. But, here we convert them into the summation term instead of the mean, and the noise covariance matrix predicted in the same way.
2. The hidden Markov transition model is hence as shown below. The extended state equation is just a simple equality, since it is some constant parameter with no evolution. We now rely on the particle filter to propose the summation term $\sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} Z_i \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h}$ instead of the mean, and the covariance matrix for V_n .

$$\hat{\mathbf{X}}_{n+1} = \begin{bmatrix} e^{\mathbf{A}\Delta t} & \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} Z_i \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} \\ 0 & 1 \end{bmatrix} \hat{\mathbf{X}}_n + \begin{bmatrix} V_n \\ 0 \end{bmatrix} \quad (37)$$

Or to be more detailed:

$$\begin{bmatrix} \mathbf{X}_{n+1} \\ \mu_w \end{bmatrix} = \begin{bmatrix} e^{\mathbf{A}\Delta t} & \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} Z_i \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \mathbf{X}_n \\ \mu_w \end{bmatrix} + \begin{bmatrix} V_n \\ 0 \end{bmatrix} \quad (38)$$

Just equation (35) but in a different convention. Note that the combined noise term would have a covariance matrix as shown in equation (39). Since the augmented state μ_w is not masked by any noise in transition.

$$\begin{bmatrix} \Sigma_{V_n} & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix} \quad (39)$$

3. Next is the emission model. That would be basically the same as before, since μ_w is not observed in this case (it is not like derivatives of \mathbf{x} , so it does not have any embedded information from it, and hence no observation?). The original observation model is:

$$\mathbf{Y}_n = \mathbf{G}\mathbf{X}_n + W_n \quad (40)$$

The new emission model is just:

$$\mathbf{Y}_n = [\mathbf{G}_n, \mathbf{0}] \begin{bmatrix} \mathbf{X}_{n+1} \\ \mu_w \end{bmatrix} + W_n \quad (41)$$

With W being the same noise as before.

4. Another important thing is the particle weight, $p(\mathbf{Y}_n | \{Z_i\}, \{\tau_i\})$. This is slightly different from the previous case, since now the jump sizes and times can no longer transform deterministically to the conditional Gaussian mean and variance. Now, they are equivalent to the summation term in the hidden Markov transition model matrix and the covariance matrix in the HMM transition noise. Therefore, conditional on the jump sizes and times, the weight probability magically turns into the marginal probability of observation again, which is perfectly compatible with our original MPF structure.

Note that this extended state technique is in fact a part of an important Kalman filtering technique called the marginalised Kalman filter.

3.3 Additional Parameter Estimation of σ_w , Margianlised Kalman Filtering

3.3.1 Problem Introduction

The main change in this step is to take out the assumption of knowing σ_w and perform a parameter estimation for it also. However, unlike μ_w which can be taken as an extended state, σ_w cannot be treated in a similar way and introduces additional uncertainty to the noise scaling in the conditional model. Although the system is still linear, but due to the assumption of known noise level in Kalman filter, we can no longer use standard Kalman filter directly to the problem. A new technique called the marginalised Kalman filter would be used here. In fact, the whole theory of taking out the μ_w parameter as an extended state and taking care of the additional uncertainty in the noise level due to σ_w are both parts of the marginalised Kalman filter technique.

Taking out the known σ_w assumption causes the particle filter to propose nothing but just the summation term, because the jump sizes and times cannot transform deterministically into the conditional Gaussian mean or variance without the NVM parameters. μ_w would be estimated in the same way as before using the same Kalman filtering scheme using an extended state. However, σ_w cannot be taken out as an extended state like μ_w , and it is some uncertain scaling factor to the covariance matrix.

$$\mathbf{S} = \sigma_W^2 \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} (\sqrt{Z_i} \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h}) (\sqrt{Z_i} \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h})^T \quad (42)$$

The main effect of this additional uncertainty is that the conditional system is no longer the usual linear model we would use with Kalman filter, and thus the weights which are observations conditional on the

particles are no longer the Kalman marginal likelihoods now, and should include the uncertainty in σ_w . This would introduce some differences in the marginal likelihood.

The solution is to work on a marginalised system that is independent of σ_w , called the marginalised system. The theory is not fully understood so far, and we would just use the structure in the Levy state space paper directly.

3.3.2 Solution Structure

We work on the same general model as before, and the main difference lies in the uncertainty distribution.

$$X_{n+1} = e^{\mathbf{A}\Delta t} X_n + \sum_{\frac{t-\Delta t}{T} \leq V_i \leq \frac{t}{T}} [\mu_w Z_i + \sigma_w \sqrt{Z_i} U_i] \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} \quad (43)$$

The particle filter basically proposes the subordinator jump sizes and times $\{Z_i\}, \{V_i\}$. The original conditional system with full knowledge of the NVM parameters would have a Gaussian noise characterised as below:

$$\mathbf{m} = \mu_w \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} Z_i \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} \quad (44)$$

$$\mathbf{S} = \sigma_w^2 \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} (\sqrt{Z_i} \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h}) (\sqrt{Z_i} \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h})^T \quad (45)$$

To add the additional uncertainty in the NVM parameters, for μ_w we take it as the additional state as before due to the clean separation of the dynamical mean, and σ_w is more tricky. We define the state vector in the conditionally linear system first:

$$\hat{\mathbf{X}} = \begin{bmatrix} \mathbf{X} \\ \mu_w \end{bmatrix} = \begin{bmatrix} x \\ \frac{dx}{dt} \\ \mu_w \end{bmatrix} \quad (46)$$

This is simply a rearrangement of terms in the original system.

$$\hat{\mathbf{X}}_{n+1} = \begin{bmatrix} e^{\mathbf{A}\Delta t} & \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} Z_i \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} \\ 0 & 1 \end{bmatrix} \hat{\mathbf{X}}_n + \begin{bmatrix} V_n \\ 0 \end{bmatrix} \quad (47)$$

Or to be more detailed:

$$\begin{bmatrix} \mathbf{X}_{n+1} \\ \mu_w \end{bmatrix} = \begin{bmatrix} e^{\mathbf{A}\Delta t} & \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} Z_i \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \mathbf{X}_n \\ \mu_w \end{bmatrix} + \begin{bmatrix} V_n \\ 0 \end{bmatrix} \quad (48)$$

The convention used in the Levy state space model paper has an additional matrix B term to transform the original noise vector V of dimension P x 1 into the new noise vector of dimension (P+1) x 1.

$$\begin{bmatrix} \mathbf{X}_{n+1} \\ \mu_w \end{bmatrix} = \begin{bmatrix} e^{\mathbf{A}\Delta t} & \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} Z_i \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \mathbf{X}_n \\ \mu_w \end{bmatrix} + B \mathbf{V}_n \quad (49)$$

$$B = \begin{bmatrix} I_{P \times P} \\ 0_{1 \times P} \end{bmatrix} \quad (50)$$

The main problem lies in the V_n term which is distributed as $N(0, \mathbf{S}_n)$ with \mathbf{S}_n defined in (45), due to the additional uncertainty in σ_w .

The marginalised Kalman filter solution defines an additional marginalised dynamical noise covariance $\bar{\mathbf{S}}_n$, such that:

$$\mathbf{S}_n = \sigma_w^2 \bar{\mathbf{S}}_n \quad (51)$$

This is just:

$$\bar{\mathbf{S}}_n = \sum_{\frac{s}{T} \leq V_i \leq \frac{t}{T}} (\sqrt{Z_i} \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h})(\sqrt{Z_i} \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h})^T \quad (52)$$

The key is that this marginalised dynamical noise term is independent of σ_w .

We define also the observation model. We follow the partial observation model mentioned in the Levy state space model paper. Here we observe y_t which is the noisy x_t with noise $W_n \sim N(0, \sigma_w^2 \kappa_v)$.

$$y_t = \mathbf{H} \hat{\mathbf{X}}_t + W_t \quad (53)$$

$$\mathbf{H} = [1, \mathbf{0}_{1 \times P}] \quad (54)$$

P is the number of original states in \mathbf{X} . Note how the observation noise W_t is assumed to scale with σ_w via some parameter κ_v . According to the Levy state space model paper, this is to ensure the greatest analytical tractability. This is important in that we also take away the assumption of knowing the observation noise level (data is still generated in the same way, but we just change to a looser observation model).

These are all the definitions needed to derive the marginalised Kalman filter.

1. First we look at the standard Kalman filter algorithm:

The transition step looks like:

$$\hat{X}_{t_i|t_{i-1}} = A \hat{X}_{t_{i-1}|t_{i-1}} \quad (55)$$

A is the standard hidden state propagation matrix in Kalman filter.

$$C_{t_i|t_{i-1}} = A C_{t_{i-1}|t_{i-1}} A^T + B S B^T \quad (56)$$

C is the hidden state uncertainty matrix inferred. S is the transition noise covariance matrix and B is the noise transformation matrix following the special convention mentioned before.

Then we look at the Kalman correction step:

$$\hat{X}_{t_i|t_i} = \hat{X}_{t_i|t_{i-1}} + K_i (y_{t_i} - H \hat{X}_{t_i|t_{i-1}}) \quad (57)$$

H is the observation matrix.

$$C_{t_i|t_i} = (I - K_i H) C_{t_i|t_{i-1}} \quad (58)$$

$$K_i = C_{t_i|t_i} H^T (H C_{t_i|t_i} H^T + C_v)^{-1} \quad (59)$$

K is the standard Kalman gain term. C_v is the observation noise covariance matrix.

2. Now, note that in our case, we have defined before $S = \sigma_w^2 \bar{S}$ and $C_v = \sigma_w^2 \kappa_v$ (**this is why that particular observation model is important. Such structure of separable σ_w^2 term is the key in analytical tractability.**) According to the paper, we can add an additional assumption that the inference covariance can also be scaled down by σ_w^2 , such that $C_{t_i|t_{i-1}} = \sigma_w^2 \bar{C}_{t_i|t_{i-1}}$. Therefore, we can cancel out the σ_w^2 term in all the covariance update steps in Kalman filter to have the marginalised Kalman filter.

In terms of the initial conditions, we would use the same ones as mentioned in the paper, with $\hat{\mathbf{X}}_{0|0} = \begin{bmatrix} \mathbf{0}_{Px1} \\ \mu_{\mu_w} \end{bmatrix}$ and $C_{0|0} = \begin{bmatrix} \mathbf{0}_{PxP} & \mathbf{0}_{Px1} \\ \mathbf{0}_{1xP} & \kappa_W \end{bmatrix}$. Note that there are two parameters κ_W and κ_v , one for the initial condition and one for the observation model.

3. Therefore, we have the marginalised Kalman filter algorithm listed as follows:

The prediction step:

$$\hat{X}_{t_i|t_{i-1}} = A\hat{X}_{t_{i-1}|t_{i-1}} \quad (60)$$

$$\bar{C}_{t_i|t_{i-1}} = A\bar{C}_{t_{i-1}|t_{i-1}}A^T + B\bar{S}B^T \quad (61)$$

The correction step:

$$\hat{X}_{t_i|t_i} = \hat{X}_{t_i|t_{i-1}} + K_i(y_{t_i} - H\hat{X}_{t_i|t_{i-1}}) \quad (62)$$

$$\bar{C}_{t_i|t_i} = (I - K_iH)\bar{C}_{t_i|t_{i-1}} \quad (63)$$

We have the new Kalman gain directly:

$$K_i = \bar{C}_{t_i|t_i}H^T(H\bar{C}_{t_i|t_i}H^T + \kappa_v)^{-1} \quad (64)$$

The observation noise covariance in this marginalised case is simply κ_v .

This is a complete Kalman filtering scheme that can be implemented directly, since the uncertain σ_w term is separated here.

The key is that, this marginalised Kalman filter scheme provides us with the accurate inference mean, but only that the uncertainty covariance matrix is σ_w^2 factor less than the exact covariance matrix. The full Kalman covariance could be restored by multiplying the marginalised covariance matrix with σ_w^2 estimate $C_{t_i|t_{i-1}} = \sigma_w^2 \bar{C}_{t_i|t_{i-1}}$. What we need now is an estimate scheme for σ_w .

3.3.3 σ_w inference framework and marginal log likelihood formula

We can hence run the σ_w estimation scheme in parallel to the marginalised Kalman filter. The method introduces a prior to the σ_w parameter using the inverted gamma distribution as shown in equation (65). Such distribution is the conjugate prior to Gaussian distribution. The reason for having a conjugate prior is that, though it can be a different distribution from the likelihood function, the combined posterior would still be of the same class as the prior distribution.

$$p(\sigma_w) = IG(\alpha_w, \beta_w) \quad (65)$$

The posterior distribution for σ_w^2 is then:

$$p(\sigma_w^2|y_{1:N}) = IG(\alpha_w + \frac{N}{2}, \beta_w + \frac{E_N}{2}) \quad (66)$$

N is the time step number. We denote $F_{t_i} = H\bar{C}_{t_i|t_{i-1}}H^T + \kappa_v$ which can be obtained directly in each step of the marginalised Kalman filter. $E_N = \sum_{i=1}^N (y_{t_i} - \hat{y}_{t_i})^T F_{t_i}^{-1} (y_{t_i} - \hat{y}_{t_i})$ which is an accumulative error term. Hence, we can compute the posterior distribution of σ_w in every time step by computing the posterior α and β .

According to the literature, the particle weights would be incremental due to the nature of inverted Gamma parameters and of the form shown in equation (67). M is the dimension of the observation vector y , N is the time step number of evaluation.

$$\log p(y_{t_1:N}) = -\frac{MN}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^N \log(|\mathbf{F}_{t_i}|) + \alpha_w \log \beta_w - (\alpha_w + \frac{N}{2}) \log \left(\beta_w + \frac{E_N}{2} \right) + \log \Gamma \left(\frac{N}{2} + \alpha_w \right) - \log \Gamma(\alpha_w) \quad (67)$$

To return the estimate of σ_w , we simply compute the mean and variance and make a plot like the Gaussian distribution. Just note that, what we are doing is inferring σ_w^2 , so there is a square 2 difference to the standard deviation we want..

Note that since such Kalman filter needs to be run for every particle, this estimation process is for every particle path, not as a whole, and it is important to use it to estimate the particle weights.

3.4 The New Hidden Markov Model with Partial Observation

Now, we change our hidden Markov model of Langevin dynamics, such that in the observation model, instead of observing the full states masked by noises, we observe one of them only, such that the other one is effectively obtained by integration or differentiation of the other state. Here, we choose to observe the state x , following the convention in the Levy state space model paper, and infer the $\frac{dx}{dt}$ state, that's essentially differentiation. The model hence looks like:

$$d\mathbf{X}(t) = \mathbf{A}\mathbf{X}(t)dt + \mathbf{h}dW(t) \quad (68)$$

The derivative relationship is contained in matrix \mathbf{A} , so there is no need to manually perform integration or differentiation, once we perform simulation or inference using the corresponding \mathbf{A} .

The solution to this SDE is again:

$$X_{n+1} = e^{\mathbf{A}\Delta t}X_n + \sum_{\frac{t-\Delta t}{T} \leq V_i \leq \frac{t}{T}} [\mu_W Z_i + \sigma_W \sqrt{Z_i} U_i] \mathbf{e}^{\mathbf{A}(-\tau_i)} \mathbf{h} \quad (69)$$

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 0 & \theta \end{bmatrix} \quad (70)$$

$$\mathbf{h} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad (71)$$

The key difference is in the observation model. Instead of the classical identity matrix which means selecting all states, it is now a row vector to have dot product with the states to select a particular state here.

$$\mathbf{Y}_n = \mathbf{G}\mathbf{X}_n + \sigma_n U_n \quad (72)$$

And now, U_n is a 1D independent Gaussian noise. \mathbf{G} is a row vector with all elements 0 except a 1 at the state selected.

$$\mathbf{G} = [1, 0] \quad (73)$$

The main change in such model is then in the Kalman correct function or other likelihood function (since we are using the marginalised particle filter here, the likelihood function is incorporated in the Kalman correction function, since the likelihood is the marginal likelihood computed in the Kalman correction step). We just change \mathbf{G} and the Cholesky decomposition of the covariance matrix of the noise, which is just the scalar std in this case.

3.5 Parameter θ estimation

Finally, we are left with the only parameter θ , the single parameter for the Langevin dynamics (could be extended directly to more complex dynamics with higher dimensional parameter space?). This would be approached via a model selection scheme by looking at the marginal probability of observations given the model parameter $p(y_{1:t}|\theta)$.

This marginal would be obtained by marginalizing the joint distribution with the state variables we are using, which are the subordinator jump sizes and times over the whole time course of data denoted as $J_{1:t}$ as shown in equation (74). It is difficult to solve such integral to give a closed form solution. However, the integral seems to fit with the particle filtering scheme and could be inferred via Monte Carlo approach.

$$p(y_{1:t}|\theta) = \int p(y_{1:t}, J_{1:t}|\theta) dJ_{1:t} = \int p(y_{1:t}|J_{1:t}, \theta) p(J_{1:t}|\theta) dJ \quad (74)$$

What the particle filter is effectively doing is proposing samples from the non-linear dynamics over each time step i.e. samples from $p(J_{t-1:t}|\theta)$. The complete particle path can be taken as a sample from $p(J_{1:t}|\theta)$. (in our case, the jump sizes and times are independent in θ , but still kept there to keep it general). The particle weights are $p(y_{t-1:t}|J_{t-1:t}, \theta)$, thus $p(y_{1:t}|J_{1:t}, \theta)$ is the total weight product. The integral can hence be interpreted as the expected total weight product. We can hence have a Monte Carlo estimate as in

equation (75), with N being the number of particles. The θ parameter inference can hence be implemented by running particle filters with different θ s.

$$p(y_{1:t}|\theta) \sim \frac{1}{N} \sum_i p(y_{1:t}|J_{1:t}^{(i)}, \theta) \quad (75)$$

4 Gaussian Distribution, Inverted Gamma Distribution and Student-t Distribution

4.1 Inverted Gamma Distribution

The Inverted Gamma (or Inverse Gamma) distribution is a continuous probability distribution that is commonly used in Bayesian statistics, especially as a conjugate prior for the variance parameter of the normal distribution. It is defined for positive real numbers and is parameterized by two parameters: the shape parameter α and the scale parameter β .

4.1.1 Probability Density Function

The probability density function (PDF) of the Inverted Gamma distribution for a random variable X is given by:

$$f(x; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{-\alpha-1} e^{-\frac{\beta}{x}}, \quad x > 0, \alpha > 0, \beta > 0$$

where $\Gamma(\alpha)$ is the Gamma function evaluated at α .

4.1.2 Mean and Variance

The mean and variance of the Inverted Gamma distribution are derived from its parameters as follows:

- **Mean:** If $\alpha > 1$,

$$\text{Mean} = \frac{\beta}{\alpha - 1}$$

- **Variance:** If $\alpha > 2$,

$$\text{Variance} = \frac{\beta^2}{(\alpha - 1)^2(\alpha - 2)}$$

Note that, since mean and variance must be positive numbers. Mean of an inverted Gamma distribution only exist for $\alpha > 1$, and variance only exist for $\alpha > 2$.

4.1.3 Conjugate Prior for the Variance of a Normal Distribution

In Bayesian statistics, the Inverted Gamma distribution is used as the conjugate prior for the variance parameter σ^2 of a normal distribution. This means that if the prior distribution of σ^2 is an Inverted Gamma distribution, then the posterior distribution of σ^2 , after observing data, is also an Inverted Gamma distribution. This property simplifies the Bayesian updating process and calculation of the posterior distribution.

When the likelihood function of the data is normal with known mean μ and unknown variance σ^2 , and the prior distribution of σ^2 is an Inverted Gamma distribution with parameters α and β , the posterior distribution of σ^2 will also be an Inverted Gamma distribution with updated parameters based on the observed data.

5 Particle MCMC

5.1 Comparison between Sequential Monte Carlo and Markov Chain Monte Carlo

Particle filtering is a **sequential Monte Carlo algorithm**, that samples from a state space system **with time evolution** like a Markov chain. **MCMC**, on the other hand, samples from a **static** distribution by using proposal to explore the sample space constrained by the acceptance probability (basically a measure of likelihood of the proposal compared to the current position). SMC samples based on the Markov structure of the state space system, whereas MCMC samples from a known distribution that is more similar to the distribution of variables at a particular time step in a Markov chain.

These are both Monte Carlo based algorithms, which implies that they can do general distributions, including those with non-linear dynamics, but they are both quite computationally expensive.

5.2 The Foundation for MCMC Algorithm: Metropolis Hasting (Reference from 4M24 Lectures)

Metropolis Hasting is an algorithm derived for sampling from the invariant distribution of a Markov chain, with the invariant distribution being the target distribution. More advanced algorithms such as Gibbs sampler and Langevin MCMC are just Metropolis Hasting using the exact proposal (acceptance probability equal to 1) and using Langevin dynamics for the proposal from a given state.

We always tend to sample from the stationary distribution of the Markov chain $\pi(dy)$. By detailed balance assuming the reversibility of the Markov Chain that:

$$\pi(x)p(x, dy) = \pi(dy)p(dy, x) \quad (76)$$

In practice, we don't always have the exact transition function $p(\cdot)$, and instead, we use a proposal function $q(\cdot)$, but the proposal transition kernel may not satisfy the equality if not perfectly match, leading to biased transition to one side. Metropolis Hasting is a solution to this problem by introducing an acceptance probability of transition $\alpha(\cdot)$:

$$\pi(x)q(x, y)\alpha(x, y) = \pi(y)q(y, x) \quad (77)$$

However, this is subject to the condition that the transition from x to y is more frequent than that from y to x for the proposal kernel $q(\cdot)$, since α is a probability and takes value from the range $[0,1]$:

$$\pi(x)q(x, y) > \pi(y)q(y, x) \quad (78)$$

In such case, we have the expression for the acceptance probability that enforces balance in the bidirectional transitions:

$$\alpha(x, y) = \frac{\pi(y)q(y, x)}{\pi(x)q(x, y)} \quad (79)$$

In practice, we don't swap x and y if the tilt is the other way round, and instead we introduce a unity cap to the acceptance probability as follows:

$$\alpha(x, y) = \min\left(\frac{\pi(y)q(y, x)}{\pi(x)q(x, y)}, 1\right) \quad (80)$$

Note that although the variables in $\pi(\cdot)$ are different in the numerator and denominator, the functions are the same and hence the normalisation constants are the same. Therefore, in computing the acceptance probability, we don't need to know the normalisation constant for π , since the one at the top and bottom cancel. (But we need for q , since these would be completely different posterior functions)

Using the acceptance probability shown in equation (80), the full proposal in MH is:

$$q_{MH}(x, y) = q(x, y)\alpha(x, y) \quad (81)$$

Note that the definition of MH still needs only the proposal, as the proposal is also the only information needed to calculate the acceptance probability. In MH, the moves to the higher probability regions are always accepted with 100% probability, but those to the lower probability regions could be rejected.

The overall algorithm for Metropolis Hasting is hence simply:

1. Simulate a sample y from $q(x^j, .)$
2. Compute $\alpha(x^j, y)$. And with the probability α (this could be obtained by generating a sample from a uniform distribution $[0,1]$ and see if the sample is greater than α), set the next sample x^{j+1} as y , otherwise set it as the same as x^j .
3. Repeat the process and end until enough samples have been obtained.

Eventually we are sampling from the invariant distribution of the Markov chain.

Gibbs sampling is just a special case of Metropolis Hasting with acceptance probability always being 1. Metropolis Hasting is the most general way to sample from a joint distribution. It also decomposes the joint distribution into a series of product of univariate conditionals which are calculated by the kernels. The only difference is that the transition could be rejected and samples may stay in the same place.

5.3 General Idea of Particle MCMC (some guess)

Particle MCMC is in essence a more efficient grid search approach. In our case, we would use it to build a more efficient parameter θ inference scheme, because as mentioned before, we can use each run of particle filter to compute a Monte Carlo expectation of the marginal likelihood $p(y_{1:t}|J^{(i)_{1:t}}, \theta)$, and this could be used to compute the acceptance probability of a new proposed θ value by taking it as the state probability $\pi(\theta)$ (there seems to be no restriction in how state i.e. θ probability should be represented).

1. We start with some initial θ parameter θ_0 . Then we run a particle filter based on this parameter to compute $p(y_{1:t}|J^{(i)_{1:t}}, \theta_0)$, i.e. $\pi(\theta_0)$
2. Then for each iteration i , we choose our MCMC algorithm to have some proposal kernel $q(., \theta_{i-1})$ to draw a θ sample, denoted as $\theta_{proposed}$. We run a particle filter for it to compute again the state probability for the proposed state $\pi(\theta_{proposed})$. The acceptance probability can be computed as in equation (82)

$$\alpha(\theta_{i-1}, \theta_{proposed}) = \min\left(\frac{\pi(\theta_{proposed})q(\theta_{proposed}, \theta_{i-1})}{\pi(\theta_{i-1})q(\theta_{i-1}, \theta_{proposed})}, 1\right) \quad (82)$$

If accepted, set θ_i to $\theta_{proposed}$, else set to θ_{i-1}

Since our distribution here is 1 dimensional only, MCMC algorithm such as Gibbs sampler would not be applicable (no lower dimensional decomposition here). We could use some simple algorithm maybe for instance, pre-conditioned Nicolson algorithm which has been used in the 4M24 coursework shown in (83). The β parameter needs to be tuned to balance the trade-off between step-size and acceptance probability.

$$p(x^{i+1}|x^i) = N(x^{i+1}; \sqrt{1 - \beta^2}x^i, \beta^2 K) \quad (83)$$