

Resampling and genealogies in sequential Monte Carlo algorithms

Susanna Elizabeth Brown

A thesis submitted for the degree of
Doctor of Philosophy in Statistics

University of Warwick, Department of Statistics

July 2021

Acknowledgements

I would like to thank...

This thesis is submitted to the University of Warwick in support of my application for the degree of Doctor of Philosophy. It has been composed by myself and has not been submitted in any previous application for any degree.

The work presented (including data generated and data analysis) was carried out by the author except in the cases outlined below:

Parts of this thesis have been published by the author:

Abstract

List of Abbreviations

SMC	sequential Monte Carlo
i.i.d.	independent and identically distributed
MRCA	most recent common ancestor
PRNG	pseudo-random number generator
CDF	cumulative distribution function
LHS	left hand side
RHS	right hand side
SSP	Srinivasan sampling process
MVB	minimal variance branching

Notation and Conventions

Also include big-O notation. And \mathbb{Z} and \mathbb{R} ?

\mathbb{N}	the natural numbers starting from one, $\{1, 2, \dots\}$
\mathbb{N}_0	the natural numbers starting from zero, $\{0, 1, 2, \dots\}$
$[a]$	the set $\{1, 2, \dots, a\}$ where $a \in \mathbb{N}$ also allow $a = 0$ in which case $[a] = \emptyset$?
\mathcal{S}_k	the k -dimensional unit simplex $\{x_{1:k+1} \geq 0 : \sum_{i=1}^{k+1} x_i = 1\}$
$(a)_b$	the falling factorial $a(a-1)\cdots(a-b+1)$ where $a \in \mathbb{N}_0, b \in \mathbb{N}$, and define $(a)_0 = 1$. could even allow $a \in \mathbb{R}$ but I don't think I ever use it in that setting
$\binom{a}{b}$	binomial coefficient where $a, b \in \mathbb{N}_0$, defined to be 0 when $a < b$
x_A	the subvector consisting of the elements of x with index in set $A \subseteq \mathbb{N}$
$x_{a:b}$	the subvector x_A where $A = \{a, a+1, \dots, b\}$
x_{-a}	the subvector x_A where $A = \{1, 2, \dots, a-1, a+1, \dots, n\}$ and n is the length of x which should be clear from the context
\prod_{\emptyset}	the empty product is taken to be 1
\sum_{\emptyset}	the empty sum is taken to be 0, while the sum over an index vector of length zero is the identity operator ?
\mathcal{F}_t	the (backward) filtration generated by offspring counts up to time t
\mathbb{E}	expectation
\mathbb{E}_t	filtered expectation $\mathbb{E}[\cdot \mid \mathcal{F}_{t-1}]$
Var	variance
Cov	covariance
A^c	the complement of set A
$ A $	the cardinality of set A
1_N	asymptotic notation for a function that converges to 1 as $N \rightarrow \infty$

2 Background

⟨ch:bg⟩

Anyone who considers arithmetical methods of producing
random digits is, of course, in a state of sin.

JOHN VON NEUMANN

2.1 Sequential Monte Carlo ✓

The idea of Monte Carlo is to use (pseudo-)random numbers to approximate expectations under an intractable probability distribution of interest. Sequential Monte Carlo (SMC) is a class of Monte Carlo algorithms which are implemented sequentially, allowing efficient sampling from *sequences* of distributions. SMC was developed for inference in intractable state space models (details in Section 2.1.1) and introduced to the statistics community by Gordon, Salmond, and Smith (1993). The basic idea behind SMC is sequential importance sampling, whereby the posterior importance samples from one target distribution are used to generate proposal samples for the next. A full derivation of the SMC recursions is beyond the scope of this work, but the reader is referred to e.g. Chopin and Papaspiliopoulos (2020) and Doucet and Johansen (2011) for more background. Here it will suffice to include a motivation in the context of state space models (Section 2.1.1) followed by the formalism of Feynman-Kac models (Section 2.1.3).

2.1.1 State space models

⟨sec:SSMs⟩

State space models (sometimes called hidden Markov models) are a flexible class of statistical models which are suitable in all sorts of applications where observations appear sequentially. The general model has two components: a Markov process $(X_t)_{t \in \mathbb{N}_0}$ representing the (unobservable) underlying state of the system, and a sequence $(Y_t)_{t \in \mathbb{N}_0}$ of noisy observations of the underlying state. The model is characterised by its conditional independence structure (Figure 2.1) along with an initial distribution μ , the Markov “transition” kernels $(K_t)_{t \in \mathbb{N}}$ and the “emission” distributions $(g_t)_{t \in \mathbb{N}_0}$. Written as a hierarchical

2 Background

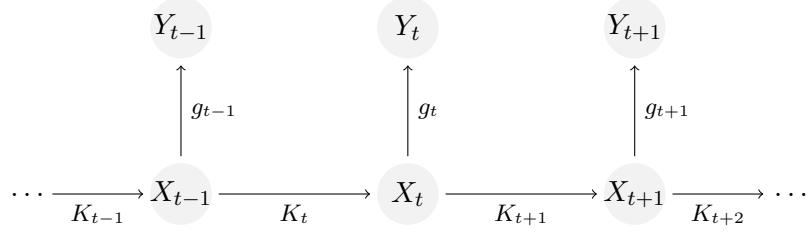


Figure 2.1: Graphical depiction of a general state space model. (X_t) is a Markov process with transition kernels (K_t) representing the underlying state of the system. Y_t is a noisy observation of X_t for each t .

(fig:SSM)

model,

$$\begin{aligned} X_0 &\sim \mu(\cdot) \\ X_{t+1} \mid X_t &\sim K_{t+1}(\cdot \mid X_t) && \text{for } t = 0, 1, \dots \\ Y_t \mid X_t &\sim g_t(\cdot \mid X_t) && \text{for } t = 0, 1, \dots \end{aligned} \tag{2.1} \quad \boxed{\text{eq:SSM_spec}}$$

The index t will frequently be referred to as “time”, since in most applications the sequence is indeed a time series, but it need not be.

Here X and/or Y may be multivariate and observation times need not be equally spaced. Straightforward generalisations of the stated model can allow for situations in which observations are not available as often as the state is updated (up to and including the extreme where the state is a continuous-time Markov process but the observations are available only at discrete times) or on the other hand where observations are observed more frequently than the state is updated.

Applications include: target tracking, where X is the true position of some object and Y encodes some measurements from sensors e.g. radar; stochastic volatility, where X is the volatility and Y is the observed value e.g. the price of a stock; change-point detection; and pretty much any other application in which there is an observed time series from which one would like to infer underlying states.

The principal inferences of interest in state space models are:

filtering $p(x_t \mid y_{0:t})$: inferring the current state x_t from the observations up to now $y_{0:t}$

prediction $p(x_{t+h} \mid y_{0:t})$: inferring a future state x_{t+h} from the observations up to now $y_{0:t}$

(complete) smoothing $p(x_{0:t} \mid y_{0:t})$: inferring the sequence of states up to now $x_{0:t}$ from the observations up to now $y_{0:t}$

fixed-lag smoothing $p(x_{t-h:t} \mid y_{0:t})$: inferring the last h states $x_{t-h:t}$ from the observations up to now $y_{0:t}$

If the dynamics of the state space model are parametrised by some θ , i.e. g_t, K_t depend on θ , we may also be interested in parameter inference and/or computing the likelihood

$p(y_{0:t})$ of the observed data given a certain value of θ .

In certain cases, these inference problems may be solved analytically (Section 2.1.2), but this is not typically the case. For intractable models we must resort to Monte Carlo methods. However, state space inference is problematic even with Monte Carlo. The main difficulties are that the dimension of the target distribution increases along the sequence, and there is strong dependence between consecutive distributions. Markov chain Monte Carlo (MCMC), for instance, is known to struggle with highly correlated targets [citation] and its performance drops drastically as dimension increases, despite convergence rates that are supposedly independent of dimension [citation].

As we will see in Section 2.1.4, sequential Monte Carlo overcomes these problems, turning the problematic properties of the target distribution to its benefit. Correlation between consecutive targets is exploited for sequential updating, which takes in its stride the incrementing dimensionality. The resulting linear-in- t computational complexity also allows inference to be performed on-line, that is, as observations arrive.

2.1.2 Exact inference in state space models

sec:SSM_exact_inference) If the state space model has linear dynamics with Gaussian errors, the posterior distributions of interest are also Gaussian with mean and covariance satisfying recursions, implemented by the Kalman filter (Kalman 1960) and Rauch-Tung-Striebel smoother (Rauch, Striebel, and Tung 1965). Recursions are also available for some other conjugate models: see for example Vidoni (1999). Another analytic case occurs if the state space \mathcal{X} is finite, in which case any integrals become finite sums, and the forward-backward algorithm (Baum et al. 1970) yields the exact posteriors. However, if the state space becomes large (albeit finite), exact computation becomes infeasible.

If the model is Gaussian but non-linear, the posterior filtering distributions can be estimated using the *extended Kalman filter* (see for example Jazwinski (2007)), which applies a first-order approximation in order to make use of the Kalman filter. This method performs well on models that are “almost linear”. The resulting predictor is only *optimal* when the model is actually linear, in which case the extended Kalman filter coincides with the Kalman filter.

For models that are high-dimensional or highly non-linear or for which gradients are not readily available, the exact Kalman filter updates can be replaced by sample approximations. The *ensemble Kalman filter* (Evensen 1994) uses a Monte Carlo sample from the current time, propagates these points through the transition dynamics, and uses the sample covariance as an estimator of the updated covariance matrix. The means (which are cheaper to evaluate and more stable than the covariances) are still updated using the Kalman filter recursion, based on the estimated covariance. The *unscented Kalman filter* (Wan and Merwe 2000) uses a deterministic sample chosen via the *unscented transformation*, which is then propagated through the non-linear transition to obtain a characterisation of the distribution at the next time step. The sample consists of $2d + 1$ points,

2 Background

where d is the dimension of the state space, and is a sufficient characterisation of a Gaussian distribution. The sample points define a Gaussian approximation to the updated distribution.

In complex or high-dimensional models, such techniques may not be feasible, in which case we must resort to Monte Carlo methods. Markov chain Monte Carlo performs woefully on state space models due to the high dimension of the parameter space and high correlation between dimensions. But we can exploit the sequential nature of the underlying dynamics to decompose the problem into a sequence of inferences of fixed dimension. This is the motivation behind sequential Monte Carlo (SMC).

2.1.3 Feynman-Kac models

(sec:FKmodels) **Example of non-SSM that is FK?**

State space models are very natural and intuitive applications, but they do not do justice to the scope of SMC algorithms, which is much wider. On the other hand, the *Feynman-Kac* formalism captures the full scope, since every SMC algorithm is a Monte Carlo approximation of some Feynman-Kac model. Before formally introducing SMC let us therefore define a generic Feynman-Kac model. For a more in-depth study, the reader is directed to the exhaustive books by Del Moral (2004, 2013) or the more accessible Chopin and Papaspiliopoulos (2020, Chapter 5).

Define a state space \mathcal{X} and a time horizon T . The basic components of the Feynman-Kac model are a Markov law, defined by an initial distribution \mathbb{M}_0 and transition kernels M_t for $t = 1, \dots, T$; and a sequence of “potential” functions $G_0 : \mathcal{X} \mapsto [0, \infty)$ and $G_t : \mathcal{X}^2 \mapsto [0, \infty)$ for $t = 1, \dots, T$. From these we construct a sequence of Feynman-Kac measures $(\mathbb{Q}_t)_{t=0:T}$ defined by the changes of measure

$$\mathbb{Q}_t(dx_{0:T}) = \frac{1}{L_t} G_0(x_0) \mathbb{M}_0(dx_0) \left\{ \prod_{s=1}^t G_s(x_{s-1}, x_s) \right\} \left\{ \prod_{s=1}^T M_s(x_{s-1}, dx_s) \right\}, \quad (2.2) \quad \text{eq:FKmeasur}$$

where L_t is the normalising constant required to make \mathbb{Q}_t a probability measure. Other quantities such as $\mathbb{Q}_t(dx_{0:t})$ can be obtained as marginals of (2.2), allowing us to treat all of the inference problems described in Section 2.1.1 by approximating \mathbb{Q}_t and then possibly marginalising.

The generic state space model introduced in (2.1) may be described by a Feynman-Kac model where:

$$\begin{aligned} \mathbb{M}_0 &:= \mu \\ M_t(x_{t-1}, dx_t) &:= K_t(dx_t \mid x_{t-1}) && \text{for } t = 1, \dots, T \\ G_0(x_0) &:= g_0(y_0 \mid x_0) \\ G_t(x_{t-1}, x_t) &:= g_t(y_t \mid x_t) && \text{for } t = 1, \dots, T. \end{aligned} \quad (2.3) \quad \text{eq:bootstra}$$

This is not the only Feynman-Kac model for (2.1); this corresponds to the “bootstrap”

2 Background

SMC algorithm, which is the simplest implementation but may be significantly outperformed by more involved algorithms such as “guided” [citation] and “auxiliary” (Carpenter, Clifford, and Fearnhead 1999; Pitt and Shephard 1999) variants. Feynman-Kac formalisms for these variants are presented for example in Chopin and Papaspiliopoulos (2020, Section 5.1.2). **Say something about the time horizon, which we have now fixed, but was infinite in SSMs section.**

It remains to demonstrate that the measures \mathbb{Q}_t arising from (2.3) are sufficient for all the usual inference problems in the corresponding state space model (2.1). By construction, the complete smoothing distribution is precisely

$$\begin{aligned}\mathbb{Q}_t(dx_{0:t}) &= \frac{1}{L_t} G_0(x_0) \mathbb{M}_0(dx_0) \prod_{s=1}^t G_s(x_{s-1}, x_s) M_s(x_{s-1}, dx_s) \\ &= g_0(y_0 \mid x_0) \mu(dx_0) \prod_{s=1}^t g_s(y_s \mid x_s) K_s(dx_s \mid x_{s-1}) \\ &= p(dx_{0:t} \mid y_{0:t}).\end{aligned}$$

The filtering, prediction and fixed-lag smoothing distributions are all also marginals of $\mathbb{Q}_t(dx_{0:T})$:

$$\begin{aligned}p(x_t \mid y_{0:t}) &= \mathbb{Q}_t(dx_t) \\ p(x_{t+h} \mid y_{0:t}) &= \mathbb{Q}_t(dx_{t+h}) \\ p(x_{t-h:t} \mid y_{0:t}) &= \mathbb{Q}_t(dx_{t-h:t}),\end{aligned}$$

while the likelihood $p(y_{0:t}) = L_t$. The upshot of this is that we have access to a Monte Carlo approximation of $\mathbb{Q}_t(dx_{0:T})$ then we can conduct inference on all of these distributions, since marginalisation of Monte Carlo samples is trivial. The likelihood, on the other hand, is not obtained by marginalisation; nevertheless, approximations of the likelihood can also be obtained “for free”. The next section describes how we may obtain samples from $\mathbb{Q}_t(dx_{0:T})$.

2.1.4 Sequential Monte Carlo for Feynman-Kac models

<sec:SMC_FK> Explain somewhere in here that the weights are supposed to give an indication of how useful each particle is, and why resampling is necessary (to avoid weight degeneracy and waste of resources on hopeless particles.)

In order to implement the SMC algorithm corresponding to a given Feynman-Kac model, we need to be able to sample from \mathbb{M}_0 and from $M_t(x, \cdot)$ for all x, t ; and evaluate $G_t(x, y)$ pointwise for each x, y, t . Under these conditions we may implement Algorithm 1, which describes a generic SMC algorithm. The only free choices are the parameter N which dictates the number of “particles” used, and the RESAMPLE procedure. However, remember that given a particular state space model there is a choice of possible Feynman-Kac

2 Background

descriptions, and this choice can strongly affect performance.

The choice of RESAMPLE procedure can also have a profound effect on performance and is discussed in detail in Section 2.4. To ensure Algorithm 1 is valid, we will assume that RESAMPLE satisfies the following properties:

- the population size N is conserved;
- given $w_{t-1}^{(1:N)}$, the expected cardinality of $\{j : a_{t-1}^{(j)} = i\}$ is proportional to $w_{t-1}^{(i)}$.

The latter property is known as *unbiasedness*. These requirements are made rigorous in Definition 2.2.

```

Input:  $T, N, \mathbb{M}_0, (M_t)_{t=1}^T, (G_t)_{t=0}^T$ 
for  $i \in \{1, \dots, N\}$  do Sample  $X_0^{(i)} \sim \mathbb{M}_0(\cdot)$ 
for  $i \in \{1, \dots, N\}$  do  $w_0^{(i)} \leftarrow \left\{ \sum_{j=1}^N G_0(X_0^{(j)}) \right\}^{-1} G_0(X_0^{(i)})$ 
for  $t \in \{1, \dots, T\}$  do
    Sample  $a_{t-1}^{(1:N)} \sim \text{RESAMPLE}(\{1, \dots, N\}, w_{t-1}^{(1:N)})$ 
    for  $i \in \{1, \dots, N\}$  do Sample  $X_t^{(i)} \sim M_t(X_{t-1}^{(a_{t-1}^{(i)})}, \cdot)$ 
    for  $i \in \{1, \dots, N\}$  do  $w_t^{(i)} \leftarrow \left\{ \sum_{j=1}^N G_t(X_{t-1}^{(a_{t-1}^{(j)})}, X_t^{(j)}) \right\}^{-1} G_t(X_{t-1}^{(a_{t-1}^{(i)})}, X_t^{(i)})$ 
end

```

Algorithm 1: Sequential Monte Carlo for a generic Feynman-Kac model

Because Algorithm 1 proceeds sequentially, its computational cost is linear in the time horizon T . Furthermore, the bootstrap algorithm, where the Feynman-Kac model is (2.3), processes the data $y_{0:T}$ one observation at a time via $G_t(x_t) = g_t(y_t | x_t)$. This means that the bootstrap SMC algorithm can also be run on-line, incorporating each observation as it becomes available. This is in stark contrast to a standard MCMC approach, for example, which would have to process all of the data at once up to a fixed time horizon. Adding one more observation would require running the MCMC algorithm from scratch on the extended target, making the computational cost (of generating samples from the whole sequence of target distributions, which is what SMC achieves) at least quadratic in time and rendering on-line inference infeasible.

The output of Algorithm 1 is, for $i = 1, \dots, N$ and $t = 0, \dots, T$, the states $X_t^{(i)} \in \mathcal{X}$ in general the state space can vary over time, the weights $w_t^{(i)} \in [0, 1]$ and, for $i = 1, \dots, N$ and $t = 0, \dots, T - 1$, the parental indices $a_t^{(i)}$. Depending on the application, one may want to retain only a subset of this output.

This output can be used to construct discrete approximations of the various distributions of interest, which can be used to estimate integrals against test functions in the usual way. The filtering distribution $p(x_t | y_{0:t})$ is approximated by

$$\sum_{i=1}^N w_t^{(i)} \delta_{X_t^{(i)}},$$

2 Background

where δ_x denotes a unit mass at x . That is, expectations of appropriate test functions $\varphi : \mathcal{X} \mapsto \mathbb{R}$ are approximated by

$$\mathbb{E}[\varphi(x_t) \mid y_{0:t}] \simeq \sum_{i=1}^N w_t^{(i)} \varphi(X_t^{(i)}).$$

The precise meaning of approximation (or \simeq) will be clarified in Section 2.1.5. To approximate the smoothing distribution, we first define the *trajectories* of states $X_{t,0:t}^{(i)}$ (for each $i \in \{1, \dots, N\}$) by setting $X_{t,t}^{(i)} := X_t^{(i)}$ and tracing back through the ancestors via the recursion $X_{t,s}(i) = X_{t,s+1}^{(a_t^{(i)})}$ for each $s \in \{0, \dots, t\}$. We then construct the approximation

$$\sum_{i=1}^N w_t^{(i)} \delta_{X_{t,0:t}^{(i)}}$$

of the smoothing distribution $p(x_{0:t} \mid y_{0:t})$. Similar approximations can be constructed for prediction and fixed-lag smoothing **...be explicit or no?**. We can also approximate the marginal likelihood using the *unnormalised* weights **should I introduce notation e.g. \tilde{w} for the unnormalised weights, and give them a separate line in Algorithm 1, and include them in the possible outputs of the algorithm?**:

$$L_t \simeq \frac{1}{N} \sum_{i=1}^N G_0(X_0^{(i)}) \prod_{t=1}^T \frac{1}{N} \sum_{i=1}^N G_t(X_{t-1}^{a_t^{(i)}}, X_t^{(i)}). \quad (2.4) \quad \boxed{\text{eq:likelihood}}$$

Define the *offspring counts* for each i, t

$$\nu_t^{(i)} := |\{j : a_t^{(j)} = i\}|,$$

the number of copies of particle i of generation t appearing in generation $t+1$ by resampling. Notice that $\nu_t^{(1:N)}$ is expressed as a non-injective function of $a_t^{(1:N)}$, and as such carries less information.

Figure 2.2 shows a section of the conditional dependence graph implied by Algorithm 1.

2.1.5 Theoretical justification

<sec:SMC_theory>

It can be shown that SMC approximations of expectations of test functions satisfy various desirable properties. For instance, it is quite easy to show that the filtering approximations converge:

$$\sum_{i=1}^N w_t^{(i)} \varphi(X_t^{(i)}) \longrightarrow \mathbb{Q}_t(\varphi),$$

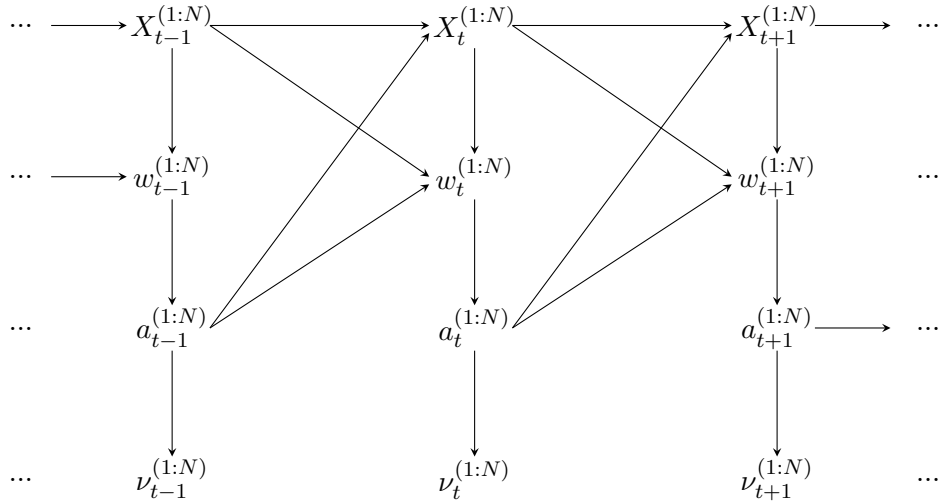


Figure 2.2: Part of the conditional dependence graph implied by Algorithm 1. The direction of time is from left to right.

`<fig:cond_indep_graph>`

almost surely and in the L_2 sense, as $N \rightarrow \infty$, under some conditions. Moreover, they satisfy a central limit theorem:

$$\sqrt{N} \left(\sum_{i=1}^N w_t^{(i)} \varphi(X_t^{(i)}) - \mathbb{Q}_t(\varphi) \right) \rightarrow \text{Normal}(0, \sigma_t(\varphi))$$

in distribution, as $N \rightarrow \infty$. Under additional conditions, the asymptotic variances $\sigma_t(\varphi)$ are stable over t , justifying the use of SMC filtering on-line. It can also easily be shown that the likelihood estimates (2.4) are unbiased (see for example Chopin and Papaspiliopoulos 2020, Proposition 16.3).

There are many other results concerning convergence, stability and error bounds for SMC algorithms. A full exposition of these results and their conditions is beyond the scope of this work, but the book by Del Moral (2013) provides an exhaustive treatment, and some of the key ideas and results are developed in the more accessible book by Chopin and Papaspiliopoulos (2020, Chapter 11). Suffice it to say that SMC algorithms enjoy enough theoretical properties to be useful in practice.

2.2 Coalescent theory ✓

`<sec:coaltheory>` Write a paragraph introducing the section.

2.2.1 Kingman's coalescent

The Kingman coalescent (Kingman 1982a,b,c) is a continuous-time Markov process on the space of partitions of \mathbb{N} . For our purposes we need only consider its restriction to $\{1, \dots, n\}$, termed the n -coalescent (defined below), since we only ever consider finite samples from a population. However, an excellent probabilistic introduction to the King-

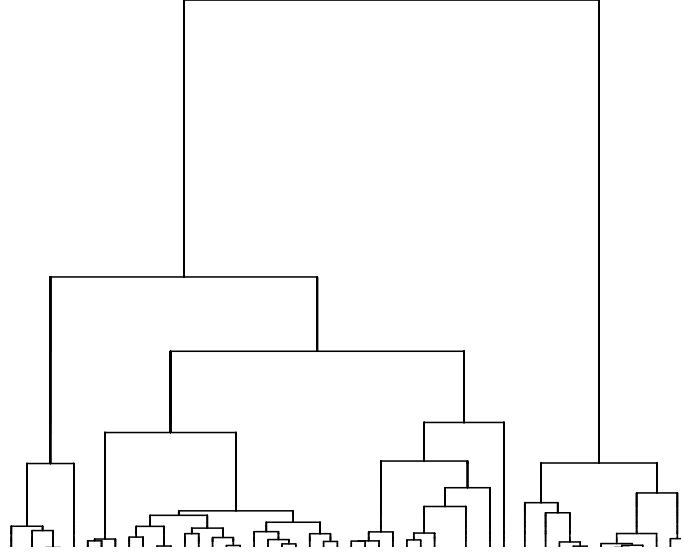


Figure 2.3: A realisation of the n -coalescent with $n = 50$.

man coalescent from the point-of-view of exchangeable random partitions can be found in Berestycki (2009, Chapters 1–2). or Wakeley (2009) ? or Durrett (2008) ?

`<def:kingman>` **Definition 2.1.** The n -coalescent is the homogeneous continuous-time Markov process on the set of partitions of $\{1, \dots, n\}$ with infinitesimal generator Q having entries

$$q_{\xi, \eta} = \begin{cases} 1 & \xi \prec \eta \\ -|\xi|(|\xi| - 1)/2 & \xi = \eta \\ 0 & \text{otherwise} \end{cases} \quad (2.5) \quad \text{?eq:KCgenerator}$$

where ξ and η are partitions of $\{1, \dots, n\}$, $|\xi|$ denotes the number of blocks in ξ , and $\xi \prec \eta$ means that η is obtained from ξ by merging exactly one pair of blocks.

A particularly attractive feature of the n -coalescent is its tractability; its distribution and those of many statistics of interest are available in closed form (Section 2.2.2). It turns out also to be extremely useful as a limiting distribution in population genetics, including the genealogies of a wide range of population models in its domain of attraction (Section 2.2.3).

2.2.2 Properties of Kingman's coalescent

`<sec:KCproperties>` Possibly also include a section about coming down from infinity (just define it basically).

The simplicity of Q allows various properties of the n -coalescent to be studied analytically. Refer to more exhaustive studies of the properties in the literature, e.g. Durrett (2008, Section 1.2). Starting with n blocks, exactly $n - 1$ coalescences are required to

2 Background

reach the absorbing state where all blocks have coalesced, known in the population genetics literature as the *most recent common ancestor* (MRCA).

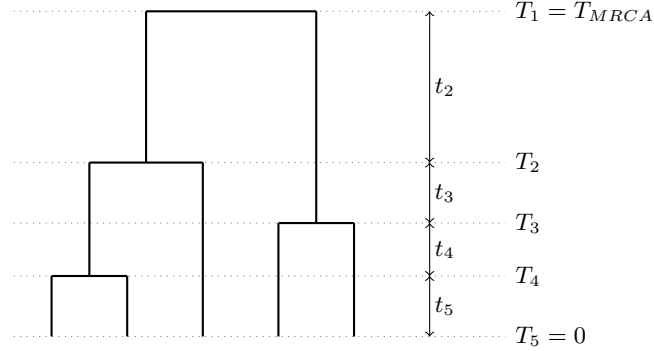


Figure 2.4: Definitions of t_i , T_i in the n -coalescent.

(fig:KC_timedefns)

Denote by t_2, t_3, \dots, t_n the waiting times between coalescent events, where t_i is the amount of time for which the coalescent has exactly i distinct lineages (see Figure 2.4). A consequence of Definition 2.1 is that these waiting times are independent and have distributions

$$t_i \sim \text{Exp} \left(\binom{i}{2} \right). \quad (2.6) \{?\}$$

The partial sum $T_k := \sum_{i=k+1}^n t_i$ gives the total time up to the $(n-k)^{\text{th}}$ coalescence event, i.e. the first time at which there are only k lineages remaining out of the initial n (see Figure 2.4). The partial sums, being sums of independent Exponential random variables, have HyperExponential distributions.

Refer back to the following three properties later on with reference to their relevance in SMC.

Time to MRCA

Of particular interest is the tree height or time to the most recent common ancestor, $T_{MRCA} := T_1$. With some algebra we find, for instance,

$$\mathbb{E}[T_{MRCA}] = \sum_{i=2}^n \mathbb{E}[t_i] = \sum_{i=2}^n \frac{2}{i(i-1)} = 2 \sum_{i=2}^n \left\{ \frac{1}{i-1} - \frac{1}{i} \right\} = 2 \left(1 - \frac{1}{n} \right) \quad (2.7) \{?\}$$

and

$$\text{Var}[T_{MRCA}] = \sum_{i=2}^n \text{Var}[t_i] = \sum_{i=2}^n \left(\frac{2}{i(i-1)} \right)^2. \quad (2.8) \{?\}$$

The expected tree height converges to 2 as $n \rightarrow \infty$, and the variance converges to $4(\pi^2 - 9)/3 \simeq 1.16$. The somewhat surprising fact that the tree height does not diverge with n is a result of the very high rate of coalescence close to the bottom of the tree. This rate is large enough that the full Kingman coalescent (on \mathbb{N}) *comes down from infinity*, that is, despite starting with infinitely many blocks, after any positive amount of time these have

2 Background

coalesced into finitely many blocks. Plot mean with sd-ribbon over n for an illustration? SD ribbon isn't the right thing; since we apparently know the actual distribution, plot a high density interval of that. (also for L)

Total branch length

Another quantity of interest is the total branch length, $L := \sum_{i=2}^n it_i$. For instance

$$\mathbb{E}[L] = \sum_{i=2}^n i \mathbb{E}[t_i] = \sum_{i=2}^n \frac{2}{i-1} = \sum_{i=1}^{n-1} \frac{2}{i} \simeq 2 \ln(n-1) \quad (2.9) \{?\}$$

and

$$\text{Var}[L] = \sum_{i=2}^n i^2 \text{Var}[t_i] = \sum_{i=2}^n \frac{4}{(i-1)^2} = \sum_{i=1}^{n-1} \frac{4}{i^2}. \quad (2.10) \{?\}$$

Note that although the mean total branch length diverges with n , the variance converges to a constant, $4\pi/6 \simeq 6.6$.

Probability that sample MRCA equals population MRCA

One other interesting quantity is the probability that the MRCA of k random lineages coincides with the population MRCA (e.g. Durrett 2008, Theorem 1.7). Denote by S_k the relevant event: that a random sample of k lineages has the same as the MRCA as the population. Consider the two subtrees produced by cutting the tree just below the population MRCA. The sample of k lineages coalesces before the population MRCA if and only if all k sampled leaves lie in just one of these two subtrees. A basic consequence of the exchangeability of the n -coalescent is that, in the limit $N \rightarrow \infty$, the proportion of leaves in the left subtree is uniformly distributed on $[0, 1]$. Calling this proportion X , we have

$$\mathbb{P}[S_k^c \mid X = x] = x^k + (1-x)^k$$

Integrating against the distribution of X , the probability of interest is

$$\mathbb{P}[S_k] = 1 - \int_0^1 [x^k + (1-x)^k] dx = \frac{k-1}{k+1}$$

as required.

The above is based on properties of the full Kingman coalescent, but similar results are available for the n -coalescent. Consider now a subsample of size k among n lineages that follow the n -coalescent. Denote by $S_{k,n}$ the event that these k lineages have the same MRCA as all n lineages. This probability of this event is calculated in Saunders, Tavaré, and Watterson (1984, Example 1) and again in Spouge (2014, Equation (3)), in both cases arising as a special case of more general results. A direct proof is given below.

Let X be the number of leaves in the left subtree. So $X \in \{1, \dots, n-1\}$ and, like before, a consequence of exchangeability is that X is uniformly distributed on that set. Now that

2 Background

the total number of branches is finite, we have to count more carefully. Conditional on X we have

$$\mathbb{P}[S_{k,n}^c \mid X = x] = \left[\binom{x}{k} + \binom{n-x}{k} \right] \binom{n}{k}^{-1}.$$

Integrating against the distribution of X gives

$$\begin{aligned} \mathbb{P}[S_{k,n}] &= 1 - \frac{1}{n-1} \binom{n}{k}^{-1} \sum_{x=1}^{n-1} \left[\binom{x}{k} + \binom{n-x}{k} \right] \\ &= 1 - \frac{1}{n-1} \binom{n}{k}^{-1} \left[\binom{n}{k+1} + \binom{n}{k+1} \right] \\ &= \frac{k-1}{k+1} \frac{n+1}{n-1} \end{aligned}$$

using binomial identities and some algebra. As $n \rightarrow \infty$ this agrees with the population-level result above.

2.2.3 Models in population genetics

`<sec:popgenmodels>`

The Kingman coalescent is the limiting coalescent process (in the large population limit) for a surprisingly wide range of population models. Some important examples of models in Kingman’s “domain of attraction” are introduced in this section. Common to all of these models are the following assumptions:

- The population has constant size N
- Reproduction happens in discrete generations
- The offspring distributions are identical at each generation, and independent between generations
- These models are all *neutral*, i.e. the offspring distribution is exchangeable.

As before [section/eq ref?](#), we define offspring counts in terms of parental indices as $\nu_j := |\{i : a_i = j\}|$. Under the assumption of neutrality, it is sufficient to consider only the offspring counts, rather than the parental indices (which generally carry more information). **Crucially, in the neutral case, offspring counts carry all the information about the distribution of the genealogy that is contained in the parental indices.** From a biological perspective, neutrality encodes the absence of natural selection, i.e. no individual in the population is “fitter” than another.

Wright-Fisher model

The neutral Wright-Fisher model (Fisher 1923, 1930; Wright 1931) is one of the most studied models in population genetics. At each time step the existing generation dies and is replaced by N offspring. The offspring descend from parents (a_1, \dots, a_N) which are

2 Background

selected according to

$$a_i \stackrel{iid}{\sim} \text{Categorical}(\{1, \dots, N\}, (1/N, \dots, 1/N)).$$

The joint distribution of the offspring counts is therefore

$$(v_1, \dots, v_N) \sim \text{Multinomial}(N, (1/N, \dots, 1/N)).$$

Since the Multinomial distribution is exchangeable, this model is neutral. There are several non-neutral variants of the Wright-Fisher model [citations?](#), but they are typically much less tractable than the neutral one.

Kingman showed in his original papers introducing the Kingman coalescent (Kingman 1982b) that, when time is scaled by a factor of N , genealogies of the neutral Wright-Fisher model converge to the Kingman coalescent as $N \rightarrow \infty$.

Cannings model

The neutral Cannings model (Cannings 1974, 1975) is a more general construction which encompasses the neutral Wright-Fisher model as a special case.

In the Cannings model, the particular offspring distribution is not specified; we only require that it is exchangeable, i.i.d. between generations, and preserves the population size. In particular, the probability of observing offspring counts (v_1, \dots, v_N) must be invariant under permutations of this vector.

Genealogies of the neutral Cannings model also converge to the Kingman coalescent, under some conditions and a suitable time-scaling [which is what?](#), as $N \rightarrow \infty$ (see for example Etheridge 2011, Section 2.2). [original reference for this? is not any Kingman 1982 papers, and certainly not Cannings 1974/5 which predates KC](#)

Moran model

The neutral Moran model (Moran 1958), while perhaps less biologically relevant, is mathematically appealing because its simple dynamics make it particularly tractable.

At each time step, an ordered pair of individuals is selected uniformly at random. The first individual in this pair dies (i.e. leaves no offspring in the next generation), while the other reproduces (leaving two offspring). All of the other individuals leave exactly one offspring. This is another special case of the neutral Cannings model, where the offspring distribution is now uniform over all permutations of $(0, 2, 1, 1, \dots, 1)$. Therefore we know that under a suitable time-scaling, its genealogies converge to the Kingman coalescent. The time scale in this case is N^2 , because reproduction happens at a rate N times [or is it technically N-1 times?](#) lower than in the Wright-Fisher model. [also cite a Moran-specific convergence result: not sure where \(it isn't in Kingman 1982* or in Moran 1958 which predates KC\)](#)

2.2.4 Particle populations

Much of the population genetics framework transfers readily to the case of SMC. The population is now a population of particles, with each iteration of the SMC algorithm corresponding to a generation, and resampling playing the part of reproduction. In fact, SMC “populations” are in some ways more suited to these population models than actual populations of organisms. The assumptions that the population has constant size N and that reproduction occurs only at discrete generations are satisfied by construction. However, we cannot assume independence between generations: as seen in Figure 2.2, the offspring counts at subsequent generations are not independent without some conditioning. In fact, after marginalising out the information about the positions of the particles, the genealogical process is not even Markovian. Nor is our model neutral: the resampling distribution depends on the weight of each particle (the weight plays the role of fitness in a non-neutral population model).

2.3 Sequential Monte Carlo genealogies ✓

`<sec:SMC_genealogies>`

We have seen that genetic terminology applies quite naturally to SMC. The resampling step induces parent-offspring relationships, each duplicate of particle i after resampling being considered one of its offspring. Then follows the notion of offspring counts (also known as family sizes), that is, the number of offspring assigned to each parent. Viewed backwards in time, the parent-offspring relationships also imply a genealogy, obtained by tracing the lineages from each terminal particle through its ancestor in each generation. We will see in this section that these genealogies, induced by resampling, are not a mere curiosity but in fact have important implications for the performance of SMC algorithms.

2.3.1 Ancestral degeneracy

Suppose we were using SMC to sample from the smoothing distribution of some state space model. As described in Section 2.1.4, we run our chosen SMC algorithm forwards, then output the N sampled trajectories $X_{t,0:t}^{(i)}$ (for each $i \in \{1, \dots, N\}$). Each trajectory was obtained by tracing back through the parent at each generation, starting from one of the terminal particles. This means that if two terminal particles i and j share a common ancestor at some generation s , then $X_{t,0:s}^{(i)}$ will be exactly equal to $X_{t,0:s}^{(j)}$, because their ancestries coincide from time s to 0.

At every resampling step, some parents may be assigned more than one offspring each, so the further back in time you look, the more of the ancestries of the terminal particles will have coalesced (see Figure 2.5a). The effect of this is that, instead of obtaining N separate sampled trajectories, we actually obtain N sampled trajectories that coalesce backwards in time, which means that the further back in time we look, the fewer distinct samples we have from the corresponding component of the target distribution. Particularly if we are interested in smoothing over a long time horizon, the variance of the SMC estimator

2 Background

is going to blow up.

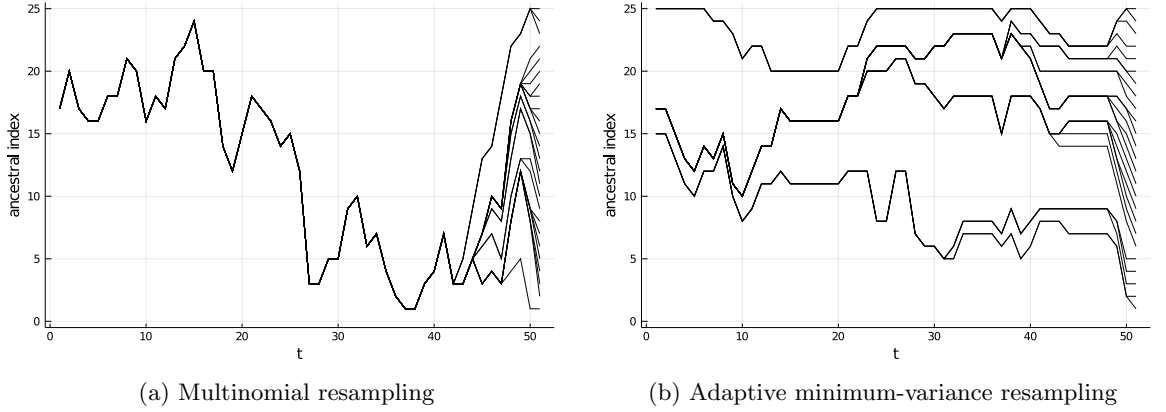


Figure 2.5: Illustration of ancestral degeneracy and the mitigating effect of low-variance and adaptive resampling. Each line is the lineage of one of the terminal particles, indicating the index of its ancestor in each generation: (a) with multinomial resampling; (b) the same system with adaptive systematic resampling.

ig:ancestral_degeneracy)

On the other hand, ancestral degeneracy actually improves the memory efficiency of SMC. We do not need to store all of the particles generated at each time (at memory cost $O(NT)$), only those that are included in the resulting genealogy. Jacob, Murray, and Rubenthaler (2015) provide an algorithm for efficient storage of the genealogy, reducing the asymptotic memory cost to $O(N \log N + T)$. However, it is certainly still worth trying to reduce ancestral degeneracy because to achieve a given level of error with a highly degenerate system will require such a large N that these memory gains are cancelled.

Mitigating ancestral degeneracy

There are a few possible approaches to mitigating ancestral degeneracy. Firstly, we could try to limit the number of offspring assigned to any one parent during each resampling step. We can only go so far, because we need the resampling procedure to remain unbiased (Section 2.1.4), but we can try to reduce the variance inherent in the resampling procedure. This idea, known as *low-variance resampling* [citation?] is discussed in detail in Section 2.4.

Another idea is to resample less often. Recall that the reason for resampling is to prevent weight degeneracy (that is, one of the weights tending to one while the others tend to zero). Now we see that, while solving one type of degeneracy, resampling creates another. The effect of ancestral degeneracy is essentially the same as that of weight degeneracy: both drastically increase the variance of the resulting SMC estimators. We can therefore consider a trade-off between the two, which is the idea behind *adaptive resampling* (Liu and Chen 1995, Section 4). The trick is to apply the resampling step only at iterations in which a certain criterion is met. The most commonly-used criterion, suggested in Liu and

2 Background

Chen (1995, Equation (14)), is based on the estimated *effective sample size*

$$ESS(t) := \left\{ \sum_{i=1}^N (w_t^{(i)})^2 \right\}^{-1},$$

which decreases as the weights degenerate. The resampling step is then applied only at iterations t such that $ESS(t)$ is less than some pre-specified threshold, typically $N/2$ [citation?].

If adaptive resampling is used, some trivial changes are required to the calculation of the weights in Algorithm 1, to allow for the importance weights to accumulate sequentially until the particles are resampled. See e.g. Chopin and Papaspiliopoulos (2020, Section 10.2) for details.

As well as mitigating ancestral degeneracy, adaptive resampling has the virtue of saving some computation (although the overall asymptotic complexity of the SMC algorithm does not change). How effective adaptive resampling is depends on the particular application and choice of SMC algorithm. If the proposals (i.e. transition kernels) are not very close to their targets then the weights will degenerate rapidly and the effective sample size criterion (or similar) will not reduce the frequency of resampling very much.

Low-variance resampling is also less effective under poor proposals: the resulting high-variance weights lead to high-variance offspring counts, even under minimum-variance resampling schemes, because the resampling is required to be unbiased.

Adaptive resampling and low-variance resampling can be combined, and this is widely considered to be the best practice when implementing SMC. Figure 2.5 compares ancestral degeneration under multinomial resampling (a relatively high variance scheme) to the same under adaptive resampling with a minimum-variance resampling scheme. It is easy to see that the degeneration is much more severe in the former case.

There is one technique that completely solves the problem of ancestral degeneracy, namely *backward simulation* [citation]. This involves running an SMC algorithm as usual (the “forward pass”), and then sampling new ancestors for each particle during an additional “backward pass”. The backward-simulated parents in each generation are chosen among all N particles, making use of particles that were not included in the forward-sampled trajectories. In terms of genealogies, the effect is striking; the lineages are broken into fragments, so there is no longer any such thing as a lineage or a genealogy.

Since this work concerns genealogies, we will not say much more about backward simulation. There are many situations in which it is impossible to implement and therefore the study of SMC genealogies is still of interest. Firstly, backward simulation inherently requires a forward and backward pass through all of the data, so it cannot be implemented on-line. Secondly, calculating the backward-simulation probabilities requires the Markov kernels M_t of the corresponding Feynman-Kac model to admit densities that can be evaluated pointwise. This is much stronger than the ability to simulate from M_t , which is the requirement for applying standard SMC algorithms.

2.3.2 Asymptotic genealogies

If we had access to information about the behaviour of SMC genealogies a priori (i.e. without having run the algorithm), we would be in a position to answer many questions of interest. These include practical questions about tuning, for example:

- How many particles should I use in order to maintain (with high enough probability) a given level of error over a time horizon T ?
- With N particles, what is the largest lag over which fixed-lag smoothing produces reasonable estimates?
- How many particles should I use within particle Gibbs to ensure that (with high enough probability) at least two distinct trajectories survive each iteration?

This last question touches on a critical aspect of the performance of particle Gibbs algorithms, which is discussed in Section 2.5. We could also consider theoretical questions, such as:

- For a given class of models and algorithms, what is the effect of ancestral degeneracy on how the estimators behave over time?
- Which resampling schemes lead to the smallest amount of ancestral degeneracy?
- What is the effect on genealogies of adaptive resampling?

Many of these questions have already been partially addressed, without any explicit analysis of genealogies, by way of variance calculations and simulation experiments. But since these are all genealogical questions by nature, it seems sensible to work directly with the genealogies, if possible. The problem is that the genealogy of particles is a complex object, it is random, and it can depend strongly on the particular choice of Feynman-Kac model and SMC implementation.

It turns out that these problems can be somewhat overcome by considering the genealogies in an asymptotic regime where the number of particles N tends to infinity. In this regime, many different particle systems exhibit genealogies of a common form, namely Kingman's n -coalescent under suitable time-scales. The differences between various algorithms is then encoded in the time-scale function, which is still random but is a less complicated object than the genealogy itself, namely a càdlàg function rather than a labelled weighted tree. In the context of SMC, these asymptotic genealogies were first analysed by Koskela et al. (2018). The simulations therein suggest that such asymptotic results also transfer to finite systems, making them practically useful.

One of the contributions of the current work is to demonstrate that Kingman-type genealogies arise from a wide variety of SMC algorithms, including those most commonly used in practice. This allows, for instance, genealogies of different SMC algorithms to be compared by examining the corresponding time-scale functions.

2.4 Resampling ~

`<sec:resampling>` As we have seen **we actually haven't; this should have been made explicit in section 2.1.4,** resampling is necessary within SMC to “reset” the weights in order to prevent weight degeneracy **and may also help us make the most of our computational resources, not wasting them on hopeless particles.** Resampling is itself a Monte Carlo procedure: the discrete offspring counts can be viewed as stochastic estimates of the continuous weights. In order to obtain a valid SMC algorithm, these Monte Carlo samples must be unbiased; this and other desirable properties are formalised in Definition 2.2. There is a huge range of resampling procedures satisfying these properties, some of which perform better than others. Some of the most popular resampling schemes are introduced in Section 2.4.2 and their properties are explored in Section 2.4.3.

2.4.1 Definition ✓

`<defn:resampling>` **Definition 2.2.** For our purposes, a valid resampling scheme is a stochastic function mapping weights $w_t^{(1:N)} \in \mathcal{S}_{N-1}$ to offspring counts $\nu_t^{(1:N)} \in \{0, \dots, N\}^N$ that satisfies the following properties:

- `item:resampling_property1>` 1. the population size is conserved: $\sum_{i=1}^N \nu_t^{(i)} = N$
- `item:resampling_property2>` 2. the weights are equal after resampling: $w_{t+}^{(i)} = 1/N$ for all i
- `item:resampling_property3>` 3. the resampling is unbiased: $\mathbb{E}[\nu_t^{(i)} \mid w_t^{(i)}] = N w_t^{(i)}$ for all i .

It is possible to design resampling schemes that violate these properties. For example, a scheme of Liu and Chen (1998) uses the square roots of the weights for resampling, then corrects by setting unequal weights after resampling (violating conditions 2 and 3). Liu, Chen, and Logvinenko (2001, Section 3.1) generalises this further, and suggests setting the resampling weights adaptively as a function of the true weights $w_t^{(1:N)}$. Fearnhead and Clifford (2003) also allows the weights to be unequal after resampling; see point (d) on p.890 therein. Resampling different (possibly random) numbers of particles in different iterations (violating condition 1) is of course possible (see for example Crisan, Del Moral, and Lyons 1999), but we typically have a fixed limit on computational resources, so in most cases it makes sense to simulate the maximum feasible number of particles N at every iteration. There may be circumstances under which it is beneficial to allow the number of particles to vary adaptively or at random, but this is not commonly done in practice. Deterministic resampling schemes (which cannot generally be unbiased, violating condition 3) have been used by some authors. These include schemes based on optimal transport (Corenflos et al. 2021; Myers et al. 2021; Reich 2013) and the importance support points resampling of Huang, Joseph, and Mak (2020). However, the majority of resampling schemes in the literature fit within Definition 2.2, and it is not typically advantageous to violate the properties 1–3.

2 Background

Within Definition 2.2 there is still a great deal of flexibility. Many different resampling schemes have been proposed in the literature, some of which perform better than others. Section 2.4.2 introduces some important resampling schemes. Their properties are discussed in Section 2.4.3 and summarised in Table 2.3.

2.4.2 Examples ✓

mples_resamplingschemes> Argue in each case that the scheme is unbiased.

Abbreviation	Description
multi	multinomial resampling
star	star resampling
strat	stratified resampling
syst	systematic resampling
res-multi	residual resampling with multinomial residuals
res-star	residual resampling with star residuals
res-strat	residual resampling with stratified residuals
res-syst	residual resampling with systematic residuals
ssp	Srinivasan sampling procedure resampling
mvp	minimal variance branching algorithm

Table 2.1: Abbreviations for resampling schemes

<tab:resampling_abbrevs>

Multinomial resampling

Multinomial resampling (Efron and Tibshirani 1994; Gordon, Salmond, and Smith 1993) is one of the simplest resampling schemes. The parental indices are chosen independently from $\{1, \dots, N\}$, each with probability given by the weight of the corresponding particle $w_t^{(i)}$. That is,

$$a_t^{(1:N)} \sim \text{Categorical}(\{1, \dots, N\}, w_t^{(1:N)}).$$

This implies that the joint distribution of the offspring counts is

$$\nu_t^{(1:N)} \stackrel{d}{=} \text{Multinomial}(N, w_t^{(1:N)}).$$

It follows from properties of the Multinomial distribution that this resampling scheme is unbiased.

A simple way to sample the parental indices is to use inversion sampling: partition the unit interval into N subintervals each of which will correspond to a certain index i and has length equal to the weight $w_t^{(i)}$; then draw N samples $U_i \sim \text{Uniform}[0, 1]$ and classify them according to which of these subintervals they fall in. Explicitly, the parental index

2 Background

assigned to child i is the index a_i satisfying

$$\sum_{j=1}^{a_i-1} w_t^{(j)} \leq U_i \leq \sum_{j=1}^{a_i} w_t^{(j)}. \quad (2.11) \quad \boxed{\text{eq:syst_str}}$$

This is illustrated in Figure 2.6.

Fast implementations of multinomial resampling rely on U_1, \dots, U_N being pre-sorted, which speeds up the search step (2.11). Sorting N numbers requires $O(N \log N)$ computation, but this is not necessary since we can sample directly the order statistics of a Uniform $[0, 1]$ distribution, at $O(N)$ cost. One way to do this (Chopin and Papaspiliopoulos 2020, Proposition 9.1) is to sample $X_i \sim \text{Exp}(1)$ independently for $i = 1, \dots, N + 1$ and output the normalised sums

$$U_k := \frac{\sum_{i=1}^k X_i}{\sum_{i=1}^{N+1} X_i}$$

for $k = 1, \dots, N$. Another method (Hol, Schön, and Gustafsson 2006) is to sample $X_i \sim \text{Uniform}[0, 1]$ for $i = 1, \dots, N$ and compute recursively

$$U_N := X_N^{1/N}, \quad U_k := X_k^{1/k} U_{k+1}.$$

I'm not sure that this really gives the correct distribution...?

This allows multinomial resampling to be implemented at $O(N)$ cost. A side-effect is that the sampled ancestral indices will be ordered and therefore cannot be Categorically distributed, but the offspring counts still have the correct Multinomial distribution. For the purposes of resampling this isn't usually a problem, but the Categorical distribution can anyway be restored at $O(N)$ cost by applying a random permutation to the offspring indices.

Residual resampling

Residual resampling is described in Liu and Chen (1998) and also in Whitley (1994) where it is called “remainder stochastic sampling”.

Each particle $X_t^{(i)}$ is deterministically assigned $\lfloor Nw_t^{(i)} \rfloor$ offspring and the remaining

$$R := \sum_{i=1}^N (Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor) = N - \sum_{i=1}^N \lfloor Nw_t^{(i)} \rfloor$$

offspring are assigned stochastically according to the residual weights

$$r^{(i)} := (Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor) / R.$$

Notice that each $r^{(i)}$ lies in the interval $[0, 1/R)$, and $R \in \{0, \dots, N - 1\}$ with $R = 0$ only if all weights are multiples of $1/N$ in which case all residual weights are zero.

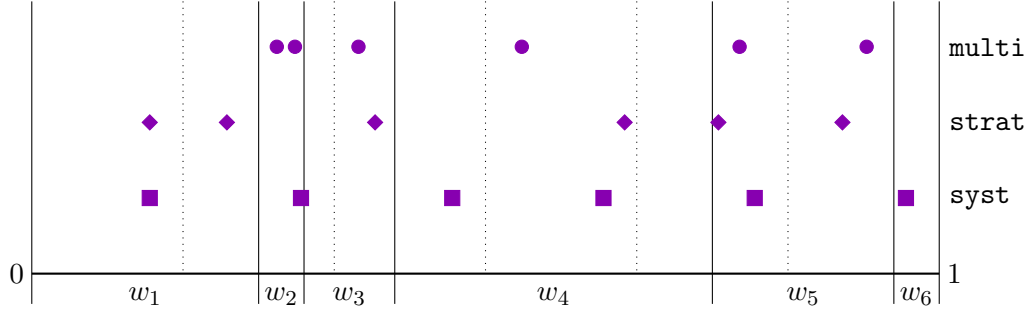


Figure 2.6: Inversion sampling interpretation of multinomial, stratified and systematic resampling. In this example, $N = 6$, $w^{(1:6)} = (0.25, 0.05, 0.1, 0.35, 0.2, 0.05)$ and the uniform random variables input to the resampling schemes are $u_{1:6} = (0.78, 0.29, 0.27, 0.92, 0.54, 0.36)$. The solid vertical lines show the partition of $[0, 1]$ into subintervals of lengths $w^{(1:6)}$. The dotted vertical lines show the partition of $[0, 1]$ into subintervals of length $1/N$, used for stratified and systematic resampling.

Top row (circles): in multinomial resampling, $u_{1:6}$ are fed directly into the inversion sampler. Which subinterval u_i falls into determines the parent of offspring i . The resulting offspring counts in this example are $\nu^{(1:6)} = (0, 2, 1, 1, 2, 0)$.

Middle row (diamonds): in stratified resampling, $u_{1:6}$ are transformed so that one point lies in each subinterval of length $1/N$. The resulting offspring counts are $\nu^{(1:6)} = (2, 0, 1, 1, 2, 0)$.

Bottom row (squares): in systematic resampling, only u_1 is used, being transformed to equally spaced points. The resulting offspring counts are $\nu^{(1:6)} = (1, 1, 0, 2, 1, 1)$.

`<fig:inv_resampling>`

The stochastic part can be implemented using any of the other basic resampling schemes (e.g. multinomial, stratified, systematic). Most presentations focus on the case where multinomial resampling is used for the residuals, which is by no means the most sensible choice. We will explore several different options in what follows.

Stratified resampling

Stratified resampling was introduced by Kitagawa (1996). As in multinomial resampling, inversion sampling is used sample the parental indices. However, the samples used for inversion sampling are no longer i.i.d. Uniform $[0, 1]$ samples. Instead, one number is sampled independently from each subinterval of length $1/N$; that is,

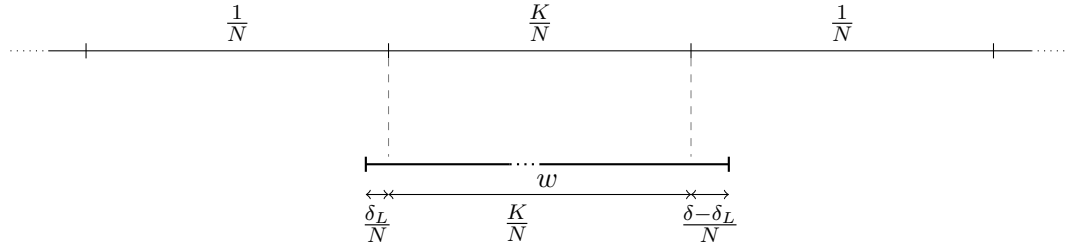
$$U_i \sim \text{Uniform}\left(\frac{i-1}{N}, \frac{i}{N}\right).$$

Alternatively, one may think of standard Uniform samples $u_1, \dots, u_N \stackrel{iid}{\sim} \text{Uniform}[0, 1]$ with the transformation

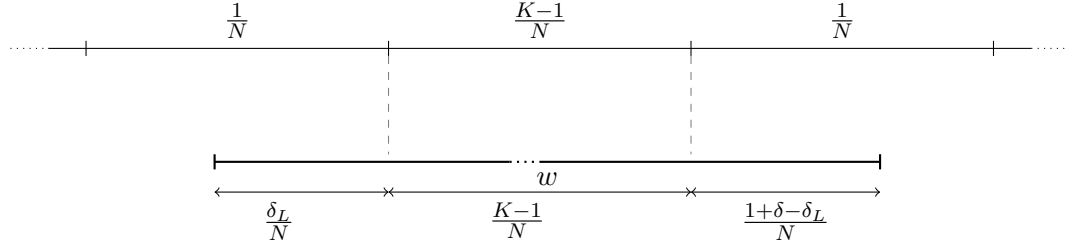
$$U_i = \frac{u_i + i - 1}{N}.$$

The parents are then assigned as in (2.11), illustrated in Figure 2.6. The offspring

2 Background



(a) The overhang is less than $1/N$ and $\delta_L \in [0, \delta]$. The parent under consideration is automatically assigned K offspring, plus up to two more.



(b) The overhang is greater than $1/N$ (this case can only occur when $K \geq 1$) and $\delta_L \in (\delta, 1)$. The parent under consideration is automatically assigned $K - 1$ offspring, plus up to two more.

Figure 2.7: Cases for stratified resampling with a fixed weight $w = (K + \delta)/N$.

`<fig:strat_cases>`

distribution is no longer Multinomial, since parental indices are not identically distributed. Stratified resampling ensures that the samples are “well spread out”, which reduces the probability of randomly losing high-weight particles or duplicating low-weight particles.

It will be useful later on to have a better idea about the marginal distributions of $\nu_t^{(i)}$ that are induced by stratified resampling. There are complex dependencies between the offspring counts, but we can still find some constraints on the distribution of each count conditional on the corresponding weight. Write the i^{th} weight in the form $w_t^{(i)} = (K + \delta)/N$, where $\delta \in [0, 1)$ and $K \in \{0, \dots, N\}$. Considering the illustration Figure 2.6, the distribution of $\nu_t^{(i)}$ depends not only on $w_t^{(i)}$ but also on where the i^{th} weight interval falls with respect to the length- $(1/N)$ intervals. Denote the “left overhang” by δ_L . There are two cases to consider, which are illustrated in Figure 2.7. In Case (a) the total overhang is less than $1/N$ and $\delta_L \in [0, \delta]$. In Case (b) the total overhang is greater than $1/N$ and $\delta_L \in (\delta, 1)$. Arrangements such that one or both ends have no overhang are special cases of Case (a) where $\delta_L \in \{0, \delta\}$. Note that Case (b) cannot occur if $K = 0$.

In any case $\nu_t^{(i)} \in \{K - 1, K, K + 1, K + 2\}$ almost surely. To define a probability distribution over these four values, we introduce the notation $p_j := \mathbb{P}[\nu_t^{(i)} = \lfloor Nw_t^{(i)} \rfloor + j \mid w_t^{(i)}]$, for $j = -1, 0, 1, 2$. Since the sample within each interval of length $1/N$ is uniform over that interval, we find the probabilities given in Table 2.2, in terms of δ and δ_L . The probabilities do not depend on K , but of course the corresponding values of $\nu_t^{(i)}$ do.

	Case (a)	Case (b)	L.B.	U.B.
p_{-1}	0	$\delta_L(1 + \delta - \delta_L) - \delta$	0	1/4
p_0	$1 - \delta + \delta_L(\delta - \delta_L)$	$1 + \delta - 2\delta_L(1 + \delta - \delta_L)$	$(1 - \delta)/2$	$1 - 3\delta/4$
p_1	$\delta - 2\delta_L(\delta - \delta_L)$	$\delta_L(1 + \delta - \delta_L)$	$\delta/2$	$(1 + \delta)/2$
p_2	$\delta_L(\delta - \delta_L)$	0	0	1/4

Table 2.2: Marginal probability distribution of $\nu_t^{(i)}$ conditional on $w_t^{(i)} = (K + \delta)/N$, in terms of δ and the left overhang δ_L , along with upper and lower bounds on these in terms of δ only, which hold in both cases.

`<tab:strat_probs>`

Systematic resampling

Systematic resampling is described in Carpenter, Clifford, and Fearnhead (1999) and also in Whitley (1994) where it is called “stochastic universal sampling”.

Like stratified resampling, it uses the inversion sampler of multinomial resampling but starts with a more regular set of points in $[0, 1]$. In this scheme, only one standard Uniform sample is drawn, $u \sim \text{Uniform}[0, 1]$, from which the N samples are generated by via the transformation

$$U_i = \frac{u + i - 1}{N}$$

for $i = 1, \dots, N$. The parental indices are again selected according to (2.11), as illustrated in Figure 2.6.

Kitagawa (1996) suggests a deterministic scheme in which the random u is replaced by a fixed $\alpha \in [0, 1]$; but, being deterministic, this scheme does not satisfy the unbiasedness property (condition 1 in Definition 2.2). Whitley (1994) employs a different description of systematic resampling, where the interval $[0, 1]$ is joined up into a circle, and the systematic samples are evenly spaced pointers on an outer ring, which is spun around like a roulette wheel. This comprises adding a random phase to each U_i , modulo one, and is an exactly equivalent description of systematic resampling.

Like stratified resampling, systematic resampling ensures the random numbers are “well spread out”; the resulting samples are even more constrained than with stratified resampling. Systematic resampling also has the advantage of being extremely easy to implement and computationally efficient, requiring only one sample from a pseudo-random number generator (PRNG) followed by $O(N)$ elementary operations.

However, the systematic scheme is known to exhibit pathological behaviour in some cases because its performance depends on the ordering of the weights. A simple example of this phenomenon is presented in Douc, Cappé, and Moulines (2005). Such behaviour can be avoided by randomly permuting the weights before resampling, and this is the recommended practice [citation?].

Star resampling

For the sake of comparison, we also construct a resampling scheme which is the worst possible (in some sense). Sample

$$a_t \sim \text{Categorical}(\{1, \dots, N\}, w_t^{(1:N)})$$

and set $a_t^{(i)} = a_t$ for all i . The resulting offspring counts are all equal to zero except for $\nu_t^{(a_t)}$, which is equal to N . This resampling scheme is indeed unbiased, since each offspring count has marginal distribution

$$\nu_t^{(i)} \mid w_t^{(1:N)} = \begin{cases} 0 & \text{w.p. } 1 - w_t^{(i)} \\ N & \text{w.p. } w_t^{(i)}. \end{cases}$$

These offspring counts have the highest possible marginal variance subject to $\mathbb{E}[\nu_t^{(i)} \mid w_t^{(i)}] = Nw_t^{(i)}$ and $\nu_t^{(i)} \in \{0, \dots, N\}$.

I call this scheme *star resampling* because the parent-offspring relationships at each iteration form a star graph.

Minimal variance branching

The minimal variance branching (MVB) algorithm of Crisan and Lyons (1999) provides a framework for resampling that enforces the minimal variance. It requires that each offspring count $\nu_t^{(i)}$, conditionally on $w_t^{(i)}$, has marginal distribution

$$\nu_t^{(i)} \mid w_t^{(i)} \stackrel{d}{=} \lfloor Nw_t^{(i)} \rfloor + \text{Bernoulli}(Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor). \quad (2.12) \quad \boxed{\text{eq:branching}}$$

We will see later on that this is exactly the framework of *stochastic rounding*.

The set-up of Crisan and Lyons (1999) does not require the number of particles to remain constant from one generation to the next (Property 1 in Definition 2.2), so the MVB algorithm could be implemented for instance by sampling each $\nu_t^{(i)}$ independently from (2.14). The authors remark that enforcing strictly negative correlation between the offspring counts can improve the rate of convergence, but they do not specify how this might be achieved. So we don't really know how to make this fit with Defn 2.2, so let's ignore it from now on?

Srinivasan sampling procedure

Gerber, Chopin, and Whiteley (2019) build on the work of Crisan and Lyons (1999) in that they construct a resampling scheme for which the marginal offspring counts are distributed as (2.14), but the number of particles is held constant and non-negative correlation of offspring counts is enforced. The resulting scheme is termed *Srinivasan sampling procedure (SSP) resampling* after Srinivasan (2001).

2 Background

The implementation is somewhat complicated compared to the other schemes we have seen (for full details see Gerber, Chopin, and Whiteley 2019, Algorithm 1) but a brief description is given here. The offspring counts are initialised at $Nw_t^{(i)}$, then we iterate through pairs of counts, rounding one of the pair up and the other down by an amount such that at least one of the pair ends up an integer. After at most N such adjustments, all of the counts are integers and can be returned. Each iteration adds and subtracts the same amount so that the sum of the counts is preserved, ensuring that the number of particles remains constant. Which of the selected pair is increased/decreased in each iteration is chosen at random with probabilities that guarantee the resampling is unbiased.

As well as proposing this resampling scheme, Gerber, Chopin, and Whiteley (2019) make several other contributions to the SMC resampling literature, some of which will be discussed later.

2.4.3 Properties ~

`c:resampling_properties` In this section we consider some important properties of resampling schemes, and see how the example schemes of Section 2.4.2 compare in terms of these. The findings are summarised in Table 2.3, with the exception of a few properties which depend on details of the implementation or are applicable only to a subset of the resampling schemes considered.

Support of offspring numbers ✓

Let us consider the support of the marginal offspring distributions in each scheme, conditional on the weights. Suppose that the i^{th} weight lies in the interval $w_t^{(i)} \in [K/N, (K + 1)/N]$.

Under multinomial resampling, it is possible for $\nu_t^{(i)}$ to take any value from 0 to N (although some values are of course more likely than others). Thus it is possible for a high-weight particle to have zero offspring, or a low-weight particle to have many offspring, simply by chance. Recall that the weights give an indication of how “useful” each particle is for the approximation. Thus killing a high-weight particle is likely to increase the variance of the SMC estimates, while duplicating a low-weight particle wastes computational resources on propagating particles that will not contribute much to reducing that variance.

Residual resampling ensures that every particle with above-average (i.e. $> 1/N$) weight has at least one offspring, avoiding the loss of high-weight particles. If the residuals are sampled using multinomial resampling then the duplication of low-weight particles is not avoided, $\nu_t^{(i)} \in \{K, \dots, K + R\} \subseteq \{K, \dots, N\}$, but this can be addressed by using a lower-variance scheme for the residual offspring. Various choices are included in Table 2.3.

Stratified resampling is more restrictive, $\nu_t^{(i)} \in \{K - 1, K, K + 1, K + 2\}$, but allows the possibility of a particle with above-average weight having no offspring. This is not quite as good as the erroneous claim of Douc, Cappé, and Moulines (2005) that $|\nu_t^{(i)} - Nw_t^{(i)}| \leq 1$ for stratified resampling. Systematic resampling has the smallest support, $\nu_t^{(i)} \in \{K, K + 1\}$,

2 Background

that is possible whilst maintaining unbiasedness, as do SSP and MVB resampling.

Another way to quantify this property is by considering the maximum possible difference between the offspring count $\nu_t^{(i)}$ and its expected value $Nw_t^{(i)}$. This is also presented in Table 2.3.

Degeneracy under equal weights ✓

In the case where all of the weights are multiples of $1/N$, low-variance schemes such as residual and systematic resampling become fully deterministic. Since $\lfloor Nw_t^{(i)} \rfloor = Nw_t^{(i)}$ for each i , residual resampling will have $R = 0$, leaving no remainder to be assigned stochastically. In systematic resampling exactly $\lfloor Nw_t^{(i)} \rfloor = Nw_t^{(i)}$ samples will fall in the i^{th} interval. In particular, if $w_t^{(1:N)} = (1, \dots, 1)/N$ then each parent is assigned exactly one offspring deterministically, so there is effectively no resampling.

The same phenomenon occurs with stratified resampling, but not if one uses Whitley’s roulette wheel description (Figure ??). The random phase shift introduced by “spinning the wheel” prevents the inversion sampling intervals from lining up exactly with the weight intervals, so the resampled offspring counts may vary from their means by one either side. Whitley (1994) does not describe stratified resampling, but we see that unlike with systematic resampling, the roulette wheel description is not equivalent to the standard inversion sampling description. For stratified resampling, the roulette wheel adds some extra randomness, so the straightforward inversion sampler is preferred.

If the state space is continuous, the event that all weights are multiples of $1/N$ typically has zero measure, but with non-zero probability we can get arbitrarily close to this regime in which resampling becomes deterministic.

Marginal variance of offspring counts ✓

Redraft this paragraph; it’s a bit clumsy atm. One indication of the performance of resampling could be the variance of the resampled offspring counts. For instance we might ask what is the marginal variance of $\nu_t^{(i)}$, conditional on the corresponding weight $w_t^{(i)}$. We want this to be low, limiting the additional randomness introduced to our Monte Carlo estimates by the resampling step.

In multinomial resampling, the marginal distributions are

$$\nu_t^{(i)} \mid w_t^{(i)} \sim \text{Binomial}(N, w_t^{(i)})$$

so the variance is

$$\text{Var}[\nu_t^{(i)} \mid w_t^{(i)}] = Nw_t^{(i)}(1 - w_t^{(i)}).$$

Compare this to star resampling, where the marginal offspring counts

$$\nu_t^{(i)} \mid w_t^{(i)} \stackrel{d}{=} N \text{Bernoulli}(w_t^{(i)})$$

2 Background

having variance

$$\text{Var}[\nu_t^{(i)} \mid w_t^{(i)}] = N^2 w_t^{(i)} (1 - w_t^{(i)}),$$

N times larger than in the multinomial case.

As pointed out in Crisan and Lyons (1999, p.557), their MVB process yields offspring variance

$$\text{Var}[\nu_t^{(i)} \mid w_t^{(i)}] = (Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor)(1 - Nw_t^{(i)} + \lfloor Nw_t^{(i)} \rfloor) \leq \frac{1}{4},$$

since the stochastic part of $\nu_t^{(i)}$ is a Bernoulli($Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor$) random variable (as seen in (2.14)). The same marginal variance appears from systematic, residual-systematic and SSP resampling, since these all share the same marginal offspring distributions. We will see in Section 2.4.4 that all of these schemes fall within the *stochastic rounding* class, and marginal offspring variance is a property shared by all stochastic roundings.

The marginal variance is harder to calculate for other schemes such as residual-multinomial and stratified resampling because these were not defined in terms of marginal distributions, nor are the offspring counts independent conditional on the weights. However, it is possible in some cases to find upper bounds on the variance, and some such bounds are derived below.

In residual-multinomial resampling, $\nu_t^{(i)}$ depends on all of the other weights as well as $w_t^{(i)}$, but only through the statistic $R := \sum (Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor)$. We have

$$\nu_t^{(i)} \mid w_t^{(i)}, R \stackrel{d}{=} \lfloor Nw_t^{(i)} \rfloor + \text{Binomial} \left(R, \frac{Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor}{R} \right).$$

Using the law of total variance,

$$\begin{aligned} \text{Var}[\nu_t^{(i)} \mid w_t^{(i)}] &= \mathbb{E} \left[\text{Var}[\nu_t^{(i)} \mid w_t^{(i)}, R] \mid w_t^{(i)} \right] + \text{Var} \left[\mathbb{E}[\nu_t^{(i)} \mid w_t^{(i)}, R] \mid w_t^{(i)} \right] \\ &= \mathbb{E} \left[(Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor) \left(1 - \frac{Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor}{R} \right) \mid w_t^{(i)} \right] \\ &\quad + \text{Var} \left[Nw_t^{(i)} \mid w_t^{(i)} \right] \\ &= Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor - (Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor)^2 \mathbb{E}[R^{-1} \mid w_t^{(i)}] \\ &\leq Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor. \end{aligned}$$

Here we have excluded the case $R = 0$, in which the variance is zero. Similarly, for residual resampling with star residuals,

$$\nu_t^{(i)} \mid w_t^{(i)}, R \stackrel{d}{=} \lfloor Nw_t^{(i)} \rfloor + R \text{Bernoulli} \left(\frac{Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor}{R} \right).$$

2 Background

and we find

$$\begin{aligned}
\text{Var}[\nu_t^{(i)} \mid w_t^{(i)}] &= \mathbb{E} \left[\text{Var}[\nu_t^{(i)} \mid w_t^{(i)}, R] \mid w_t^{(i)} \right] + \text{Var} \left[\mathbb{E}[\nu_t^{(i)} \mid w_t^{(i)}, R] \mid w_t^{(i)} \right] \\
&= \mathbb{E} \left[R(Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor) \left(1 - \frac{Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor}{R} \right) \mid w_t^{(i)} \right] \\
&\quad + \text{Var} \left[Nw_t^{(i)} \mid w_t^{(i)} \right] \\
&= \mathbb{E} \left[R(Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor) \left(1 - \frac{Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor}{R} \right) \mid w_t^{(i)} \right] \\
&= (Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor) \mathbb{E} \left[R \mid w_t^{(i)} \right] - (Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor)^2 \\
&\leq N(Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor).
\end{aligned}$$

Again, if $R = 0$ then the variance is zero.

For stratified resampling, we can use the constraints on the marginal offspring distribution that were derived in Section 2.4.2. Recall that, conditional on $w_t^{(i)}$, $\nu_t^{(i)} = \lfloor Nw_t^{(i)} \rfloor + j$ with probability p_j for $j = -1, 0, 1, 2$. We can use the expressions for p_{-1}, p_0, p_1, p_2 in the two cases of Figure 2.7, as summarised in Table 2.2, to bound the variance. First write

$$\begin{aligned}
\text{Var} \left[\nu_t^{(i)} \mid w_t^{(i)} \right] &= \mathbb{E} \left[(\nu_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor)^2 \mid w_t^{(i)} \right] - \mathbb{E} \left[\nu_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor \mid w_t^{(i)} \right]^2 \\
&= p_{-1} + p_1 + 4p_2 - (-p_{-1} + p_1 + 2p_2)^2.
\end{aligned} \tag{2.13} \quad \boxed{\text{eq:marg_var}}$$

Using the upper and lower bounds in Table 2.2 and then optimising over δ , we obtain the bound

$$\text{Var}[\nu_t^{(i)} \mid w_t^{(i)}] \leq \frac{1}{4} + \frac{1+\delta}{2} + 1 - (0 + \frac{\delta}{2} + 0)^2 = \frac{1}{4}(7 + 2\delta - \delta^2) \leq 2.$$

Optimising the exact expressions in each case (first two columns in Table 2.2) does not improve this overall bound.

Residual-stratified resampling has the further constraint that $p_{-1} = 0$ (i.e. Figure 2.7b doesn't occur) since the residual weights are between 0 and $1/R$. Now the bounds in Table 2.2 are too loose, so we bound the variance by using the exact expressions from Table 2.2 in each case and optimising over δ_L, δ . By setting p_{-1} to zero in (2.15), substituting the expressions for Case (a) from Table 2.2, and maximising over δ_L and then δ , we get

$$\begin{aligned}
\text{Var}[\nu_t^{(i)} \mid w_t^{(i)}] &= p_1 + 4p_2 - (p_1 + 2p_2)^2 = \delta - 2\delta_L(\delta - \delta_L) + 4\delta_L(\delta - \delta_L) - \delta^2 \\
&= \delta - \delta^2 + 2\delta\delta_L - 2\delta_L^2 \leq \delta - \frac{1}{2}\delta^2 \leq \frac{1}{2}.
\end{aligned}$$

Similarly, for Case (b),

$$\text{Var}[\nu_t^{(i)} \mid w_t^{(i)}] = p_1 + 4p_2 - (p_1 + 2p_2)^2 = \delta_L(1 + \delta - \delta_L) - (\delta_L(1 + \delta - \delta_L))^2 \leq \frac{1}{4}.$$

2 Background

Combining the two cases, we obtain an overall variance bound

$$\text{Var}[\nu_t^{(i)} \mid w_t^{(i)}] \leq \frac{1}{2}$$

for residual-stratified resampling.

Table 2.3 includes upper bounds on $\text{Var}[\nu_t^{(i)}]$ for various resampling schemes, independent of $w_t^{(i)}$. Those general bounds are derived from the results of this section, bounded above independently of the weights. Some of the bounds may not be tight. We could also try to bound this variance below, but for every resampling scheme the only lower bound valid for all $w_t^{(i)}$ is zero (consider the case $w_t^{(i)} = 0$) so this would not provide any more information.

Contribution to the Monte Carlo variance ✓

While the variance of the offspring counts goes some way towards providing a comparison between the various resampling schemes, a more relevant property is the contribution of the resampling step to the Monte Carlo variance. This quantifies directly the effect of a certain choice of resampling scheme on the variance of the resulting Monte Carlo estimators.

Let $(\mathcal{G}_t)_{t \geq 0}$ be the filtration generated by the particle positions and weights up to and including time t , so \mathcal{G}_t is the σ -algebra generated by $(X_{0:t}^{(1:N)}, w_{0:t}^{(1:N)})$. Consider the position of the i th particle in generation $t + 1$ just after resampling but before mutating, that is $X_t^{(a_t^{(i)})}$. Define the one-step Monte Carlo variance induced by resampling as

$$\sigma(\varphi) := \text{Var} \left[\frac{1}{N} \sum_{i=1}^N \varphi(X_t^{(a_t^{(i)})}) \mid \mathcal{G}_t \right] \quad (2.14) \quad \boxed{\text{eq:resampling}}$$

where φ is an arbitrary test function.

Some results comparing this variance across different resampling schemes are presented in Douc, Cappé, and Moulines (2005). Their results, plus some additional ones, are presented in Proposition 2.3. It may be possible to derive similar results regarding residual-stratified and SSP resampling, but such results are hard to obtain due to the strong dependence between parental indices induced by these resampling schemes. This remains an interesting open problem.

In the case of systematic (but not necessarily residual-systematic) resampling, no such variance comparison can be made. Systematic resampling generally yields low variance in practice, but it is possible to construct pathological cases in which it yields higher variance than multinomial resampling (Douc, Cappé, and Moulines 2005, Section 3.4) and it lacks theoretical support more generally (e.g. Gerber, Chopin, and Whiteley 2019, Section 3.3).

2 Background

hm:resampling_var_compare)

Proposition 2.3 (Variance of resampling schemes). *Let σ_{multi} etc. denote the variance (2.16) under the various resampling schemes, as abbreviated in Table 2.1. For any square-integrable function φ ,*

$\langle \text{item:resampling_var1} \rangle$

$$(a) \quad \sigma_{\text{multi}}(\varphi) \geq \sigma_{\text{res-multi}}(\varphi)$$

$\langle \text{item:resampling_var2} \rangle$

$$(b) \quad \sigma_{\text{multi}}(\varphi) \geq \sigma_{\text{strat}}(\varphi)$$

$\langle \text{item:resampling_var3} \rangle$

$$(c) \quad \sigma_{\text{star}}(\varphi) = N\sigma_{\text{multi}}(\varphi)$$

$\langle \text{item:resampling_var4} \rangle$

$$(d) \quad \sigma_{\text{res-star}}(\varphi) \geq \sigma_{\text{res-multi}}(\varphi) \geq \sigma_{\text{res-strat}}(\varphi)$$

Proof. (a) See Douc, Cappé, and Moulines (2005, Section 3).

(b) See Douc, Cappé, and Moulines (2005, Section 3).

(c) The following expression is derived in Douc, Cappé, and Moulines (2005, Equation (6)):

$$\sigma_{\text{multi}}(\varphi) = \frac{1}{N} \sum_{j=1}^N \varphi^2(X_t^{(j)}) w_t^{(j)} - \frac{1}{N} \left\{ \sum_{j=1}^N \varphi(X_t^{(j)}) w_t^{(j)} \right\}^2.$$

Under star resampling, all of the resampled indices are equal, say $X_t^{(a_t^{(1)})} = \dots = X_t^{(a_t^{(N)})} = X_t^*$, so

$$\begin{aligned} \sigma_{\text{star}}(\varphi) &= \text{Var} \left[\frac{1}{N} \sum_{i=1}^N \varphi(X_t^{(a_t^{(i)})}) \mid \mathcal{G}_t \right] = \text{Var} [\varphi(X_t^*) \mid \mathcal{G}_t] \\ &= \mathbb{E} [\varphi^2(X_t^*) \mid \mathcal{G}_t] - \mathbb{E} [\varphi(X_t^*) \mid \mathcal{G}_t]^2 \\ &= \sum_{j=1}^N \varphi^2(X_t^{(j)}) \mathbb{P}[X_t^* = X_t^{(j)} \mid \mathcal{G}_t] - \left\{ \sum_{j=1}^N \varphi(X_t^{(j)}) \mathbb{P}[X_t^* = X_t^{(j)} \mid \mathcal{G}_t] \right\}^2 \\ &= \sum_{j=1}^N \varphi^2(X_t^{(j)}) w_t^{(j)} - \left\{ \sum_{j=1}^N \varphi(X_t^{(j)}) w_t^{(j)} \right\}^2 \\ &= N\sigma_{\text{multi}}(\varphi), \end{aligned}$$

as required.

(d) The second inequality follows from (b) and is stated in Gerber, Chopin, and Whiteley (2019, p.9). For the first inequality, we use the following expression which is a slight modification of Douc, Cappé, and Moulines (2005, Equation (8)):

$$\sigma_{\text{res-multi}}(\varphi) = \frac{R}{N^2} \sum_{j=1}^N \varphi^2(X_t^{(j)}) r^{(j)} - \frac{R}{N^2} \left(\sum_{j=1}^N \varphi(X_t^{(j)}) r^{(j)} \right)^2.$$

A derivation similar to theirs can also be used for residual-star resampling. First notice

2 Background

that, conditional on \mathcal{G}_t , the Monte Carlo estimate in (2.16) can be decomposed into a sum of conditionally deterministic terms plus a sum of stochastic terms:

$$\frac{1}{N} \sum_{i=1}^N \varphi(X_t^{(a_t^{(i)})}) = \frac{1}{N} \sum_{j=1}^N \lfloor N w_t^{(j)} \rfloor \varphi(X_t^{(j)}) + \frac{1}{N} \sum_{i=1}^R \varphi(\hat{X}_t^{(i)}),$$

where the terms in the second sum are all equal, say $\hat{X}_t^{(1)} = \dots = \hat{X}_t^{(R)} = X_t^*$. The first sum is conditionally deterministic and hence does not contribute to the Monte Carlo variance (2.16). We have

$$\begin{aligned} \sigma_{\text{res-star}}(\varphi) &= \text{Var} \left[\frac{1}{N} \sum_{i=1}^R \varphi(\hat{X}_t^{(i)}) \middle| \mathcal{G}_t \right] = \frac{R^2}{N^2} \text{Var} [\varphi(X_t^*) | \mathcal{G}_t] \\ &= \frac{R^2}{N^2} \mathbb{E} [\varphi^2(X_t^*) | \mathcal{G}_t] - \frac{R^2}{N^2} \mathbb{E} [\varphi(X_t^*) | \mathcal{G}_t]^2 \\ &= \frac{R^2}{N^2} \sum_{j=1}^N \varphi^2(X_t^{(j)}) \mathbb{P}[X_t^* = X_t^{(j)} | \mathcal{G}_t] - \frac{R^2}{N^2} \left\{ \sum_{j=1}^N \varphi(X_t^{(j)}) \mathbb{P}[X_t^* = X_t^{(j)} | \mathcal{G}_t] \right\}^2 \\ &= \frac{R^2}{N^2} \sum_{j=1}^N \varphi^2(X_t^{(j)}) r^{(j)} - \frac{R^2}{N^2} \left\{ \sum_{j=1}^N \varphi(X_t^{(j)}) r^{(j)} \right\}^2 \\ &= R \sigma_{\text{res-multi}}(\varphi) \\ &\geq \sigma_{\text{res-multi}}(\varphi) \end{aligned}$$

whenever $R \geq 1$. If $R = 0$ then all residual schemes have zero variance and (d) holds trivially. ■

Exchangeability of offspring ✓

We say that a resampling scheme leaves the offspring exchangeable if the resulting distribution of parental indices is invariant under permutations of the offspring. To put it another way, each child chooses its parent from the same marginal distribution.

It is clear that true multinomial resampling satisfies this property since the parental indices are independent and distributed according to the same Categorical distribution. The same goes for star resampling. However, as mentioned earlier, the efficient implementation of multinomial resampling that takes sorted inputs does not leave the offspring exchangeable. Stratified and systematic resampling do not either since their inversion sampling points are sorted: for instance, child 1 is more likely to choose parent 1 than child N is. Residual resampling schemes are also typically implemented in such a way that the offspring are not exchangeable.

Whichever resampling scheme is used, exchangeability of offspring can easily be reintroduced, at $O(N)$ cost, by applying a random permutation to the vector of parental indices after sampling.

2 Background

Operations in SMC that depend on the ancestral indices are typically independent of ordering, so sampling ancestral indices from a non-exchangeable distribution is not expected to cause any problem. However, the results of Chapters 3 and 4 rely on assumption (A1) which amounts to exchangeability of offspring, so to be sure that the current genealogical study applies, a permutation should be appended to any non-exchangeable resampling procedure.

Permutation sensitivity ~

[link back to systematic](#), where we said that permutation sensitivity can cause pathological behaviour. Some resampling schemes are sensitive to the order of the weights. That is, permuting the weight vector before resampling can affect the distribution of the resulting offspring counts. Note that this entails a permutation of the parents, whereas the previous section was about permutations of offspring.

For example, consider resampling schemes based on inversion sampling (multinomial, stratified, systematic). Figure 2.8 shows two partitions of $[0, 1]$ each constructed from a permutation of the weight vector $w^{(1:6)} = (0.25, 0.05, 0.1, 0.35, 0.2, 0.05)$. This does not affect the distribution of the offspring counts under multinomial resampling, although it will affect the distribution of the parental indices if the fast implementation is used.

On the other hand, under stratified or systematic resampling the distribution of offspring counts is different for the two partitions. To see this, consider parents 2 and 6. When the weights are sorted, the probability that both of these parents are assigned a non-zero number of offspring is zero, because both of their subintervals lie within the same subinterval of length $1/N$, which gets exactly one inversion sampling point. When the weights are in their natural order, as in the top row of Figure 2.8, it is possible under stratified and systematic resampling for both parents 2 and 6 to be assigned one offspring. Clearly, then, the distribution of offspring counts under these resampling schemes differs between the two orderings of the weight vector pictured.

Which other schemes do/don't have this property?

What effect does permutation sensitivity have on performance? — perhaps we just lead into the “Sortin” section to answer that...

Sorting

Results from Gerber, Chopin, and Whiteley (2019) about benefits of sorting. What about sorting instead by weights?

Computational complexity ✓

All of the resampling algorithms discussed in Section 2.4.2 can be implemented in $O(N)$ operations. Considering the complexity of each operation, Hol, Schön, and Gustafsson (2006) suggest that systematic resampling is fastest because it only requires one pseudo-random number generation, and multinomial resampling is slower than stratified resam-

2 Background

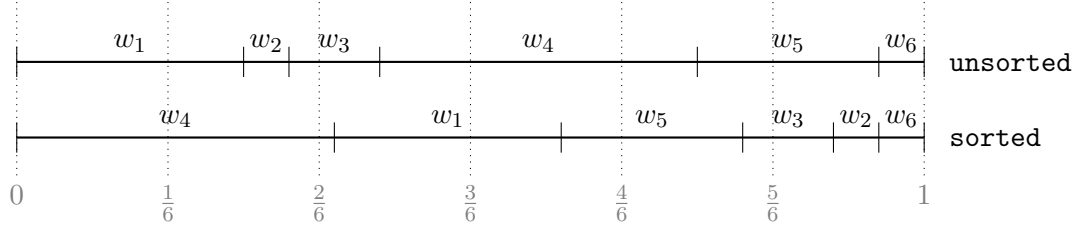


Figure 2.8: An example in which permuting the weights can affect the conditional distribution of offspring counts under certain resampling schemes. As in Figure 2.6, $N = 6$ and $w^{(1:6)} = (0.25, 0.05, 0.1, 0.35, 0.2, 0.05)$. The top row shows the weighted subintervals in the natural order, as in Figure 2.6. The bottom row shows the partition corresponding to the same weights, but sorted in decreasing order. The dotted lines are spaced $1/N$ apart. Under “permutation-sensitive” resampling schemes, the distribution of offspring counts differs depending on which partition is used.

permutation_sensitivity)

pling because of the transformations required (although this may depend on which fast implementation of multinomial sampling is used). Residual resampling is hard to compare directly because a random fraction of the operations are deterministic, so the number of pseudo-random numbers required is a random number between 0 and $N - 1$, but the authors’ simulation experiments place it between multinomial and stratified resampling.

However, the analysis of per-particle cost is sensitive to the particular implementation of each resampling scheme, the system implementation of pseudo-random number generation and arithmetic operations, and the hardware used, so it is not clear how robust such comparisons are.

Negative association ✓

Following Gerber, Chopin, and Whiteley (2019), we use the definition of negative association from Joag-Dev and Proschan (1983).

Definition 2.4. Let (Z_1, \dots, Z_n) be a collection of random variables. $Z_{1:n}$ are said to be *negatively associated* if, for every disjoint pair of subsets $I, J \subseteq \{1, \dots, n\}$, for all real-valued coordinatewise non-decreasing functions φ, ψ for which the covariance is well defined,

$$\text{Cov}[\varphi(Z_I), \psi(Z_J)] \leq 0.$$

Gerber, Chopin, and Whiteley (2019) show that negative association of offspring counts is a desirable property which may be used, along with some other machinery, to establish certain weak convergence results for the resampled measures.

Multinomial counts are negatively associated (Joag-Dev and Proschan 1983, Section 3.1), which implies that residual-multinomial resampling also satisfies this property. **Why is that? check it’s true! idea: deterministic parts have 0 correlation and the random parts are NA. same question for the other residual schemes mentioned below.** Gerber,

2 Background

Chopin, and Whiteley (2019) construct a counter-example to demonstrate that systematic resampling violates the negative association property. For residual-systematic resampling, we can cook up a counterexample in the same spirit by taking $\varphi(x) = \psi(x) = \mathbb{1}_{\{x=1\}}$, $I = \{1\}$, $J = \{3\}$ and considering a weight vector say $w^{(1:4)} = \frac{1}{8}(1, 1, 1, 5)$ for $N = 4$. Then the residual weights are $r^{(1:4)} = \frac{1}{4}(1, 1, 1, 1)$ with $R = 2$, so

$$\begin{aligned} \text{Cov}[\varphi(Z_I), \psi(Z_J)] &= \mathbb{E}[\varphi(Z_I)\psi(Z_J)] - \mathbb{E}[\varphi(Z_I)]\mathbb{E}[\psi(Z_J)] \\ &= \mathbb{P}[\nu^{(1)} = 1, \nu^{(3)} = 1] - \mathbb{P}[\nu^{(1)} = 1]\mathbb{P}[\nu^{(3)} = 1] \\ &= \frac{1}{2} - \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} > 0, \end{aligned}$$

since $\nu^{(3)} = 1$ if and only if $\nu^{(1)} = 1$. So residual-systematic resampling also violates the negative association property.

Gerber, Chopin, and Whiteley (2019) also mention some resampling schemes that do result in negatively associated counts: stratified resampling, and by implication residual-stratified resampling; star resampling (see the remark at the end of Gerber, Chopin, and Whiteley (2019, Section 3.2)), and by implication residual-star resampling. The authors go on to introduce the SSP resampling scheme, which yields negatively associated offspring counts by construction.

These results are summarised in Table 2.3. The MVB algorithm does not enforce negative association, so this property depends on the particular implementation, and as such is left blank in Table 2.3.

Star discrepancy ✓

The *star discrepancy* is a measure of the regularity of a given set of points $u_{1:N}$ in the unit hypercube. For our purposes it is sufficient to define the star discrepancy in one dimension, as in Kuipers and Niederreiter (1974, Definition 1.2):

$$D^*(u_1, \dots, u_N) := \sup_{u \in [0,1]} |d(u)| := \sup_{u \in [0,1]} \left| u - \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{u_i \leq u\}} \right|. \quad (2.15) \quad \boxed{\text{eq: defn_star}}$$

The quantity inside the supremum is the difference between the empirical CDF of the observed points $u_{1:N}$ and the CDF of the Uniform distribution on $[0, 1]$. Thus D^* measures, in a certain sense, how far the points are from being uniformly spaced.

Star discrepancy is used in quasi-Monte Carlo, where “low-discrepancy” points are used in place of uniform samples to decrease the variance of Monte Carlo estimates. We have noted already that resampling can itself be viewed as a Monte Carlo procedure. From this point-of-view, stratified and systematic resampling are quasi-Monte Carlo versions of multinomial resampling, since they provide “more regular” point sets to be used in inversion sampling.

In one dimension, the lowest-discrepancy point set is the regular grid $\frac{1}{2N}(1, 3, \dots, 2N - 1)$, which has star discrepancy $\frac{1}{2N}$ (see for example Kuipers and Niederreiter 1974, Corol-

2 Background

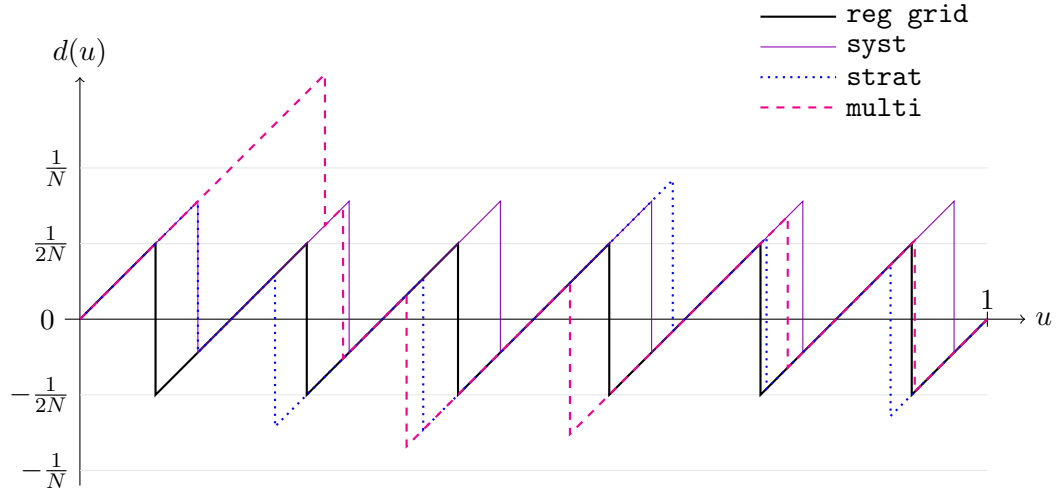


Figure 2.9: Plot of the function inside the absolute value in (2.17), for four different point sets. The points $u_{1:6}$ used are the same as in Figure 2.6.

The solid black line corresponds to the regular grid, which achieves the minimal discrepancy $1/(2N)$, but cannot be used for resampling. The star discrepancy of stratified and systematic points varies between $1/(2N)$ and $1/N$ depending on the realisation. In this example, the star discrepancy of the systematic points is $0.78/N$ and of the stratified points is $0.92/N$. The star discrepancy of standard multinomial resampling (that is, i.i.d. Uniform points) can be arbitrarily close to 1 for “bad” realisations; in this example it is $1.62/N$.

⟨fig:star_discrepancy⟩

lary 1.2). However, resampling based on a deterministic point set cannot be unbiased since the resulting parental indices are conditionally deterministic given the weights. Systematic resampling amounts to a randomisation of the regular grid, shifting each grid point by a random amount $u \sim \text{Uniform}[0, 1/N]$. This yields star discrepancy $D^* = \max\{u, \frac{1}{N} - u\}$, which is between $1/(2N)$ and $1/N$ almost surely. The point sets generated in stratified resampling also have star discrepancy between $1/(2N)$ and $1/N$, where the exact value depends on the realisation. This certainly seems to improve on independent uniform points which can have star discrepancy arbitrarily close to 1, the maximum possible value, albeit with diminishing probability as N increases. Figure 2.9 illustrates how the star discrepancy is computed, and how it compares between these sampling methods.

	support of $\nu_t^{(i)}$ given $\frac{K}{N} \leq w_t^{(i)} < \frac{K+1}{N}$	$\sup_w \nu_t^{(i)} - Nw_t^{(i)} $	upper bound on $\text{Var}[\nu_t^{(i)}]$	stochastic rounding?	degenerate if $w_t^{(1:N)} = \frac{1}{N}(1, \dots, 1)$?	sensitive to permutations of weights?	PRNG calls	neg. assoc.?
multi	$\{0, \dots, N\}$	N	$N/4$	\times	\times	\times	N	\checkmark
star	$\{0, N\}$	N	$N^2/4$	\times	\times	\times	1	$\checkmark?$
strat	$\{K-1, K, K+1, K+2\}$	2	2	\times	\checkmark	\checkmark	N	\checkmark
syst	$\{K, K+1\}$	1	$1/4$	\checkmark	\checkmark	\checkmark	1	\times
res-multi	$\{K, \dots, N\}$	N	1	\times	\checkmark	\times	$\leq N-1$	\checkmark
res-star	$\{K, N\}$	N	N	\times	\checkmark	\times	1	$\checkmark?$
res-strat	$\{K, K+1, K+2\}$	2	$1/2$	\times	\checkmark	\checkmark	$\leq N-1$	\checkmark
res-syst	$\{K, K+1\}$	1	$1/4$	\checkmark	\checkmark	\checkmark	1	\times
ssp	$\{K, K+1\}$	1	$1/4$	\checkmark	\checkmark	$\checkmark?$?	\checkmark
mvb	$\{K, K+1\}$	1	$1/4$	\checkmark	\checkmark			

Table 2.3: Summary of some of the properties of resampling schemes explored in Section 2.4.3. The abbreviated names for the resampling schemes are explained in Table 2.1. **I need to include an explanation of the column titles in the caption too.** Some properties are not specified for branching because they depend on the particular implementation.

2.4.4 Stochastic rounding ✓

⟨sec:SRs⟩

?⟨defn:stochround⟩?

Definition 2.5. Let $X = (X_1, \dots, X_N)$ be a \mathbb{R}_+^N -valued random variable. Then $Y = (Y_1, \dots, Y_N) \in \mathbb{N}^N$ is a *stochastic rounding* of X if each element Y_i takes values

$$Y_i \mid X_i = \begin{cases} \lfloor X_i \rfloor & \text{with probability } 1 - X_i + \lfloor X_i \rfloor \\ \lfloor X_i \rfloor + 1 & \text{with probability } X_i - \lfloor X_i \rfloor. \end{cases}$$

By construction, $\mathbb{E}(Y_i) = X_i$ for each i . Taking X to be N times the vector of particle weights, we can therefore use stochastic rounding to construct a valid resampling scheme, under the further constraint that $Y_1 + \dots + Y_N = N$. Several ways to enforce this constraint on the joint distribution have been proposed, including systematic resampling, residual resampling with systematic residuals, and SSP resampling.

Explicitly, the offspring counts are marginally distributed according to

$$\nu_t^{(i)} \mid w_t^{(i)} \stackrel{d}{=} \lfloor Nw_t^{(i)} \rfloor + \text{Bernoulli}(Nw_t^{(i)} - \lfloor Nw_t^{(i)} \rfloor).$$

Some of the properties discussed earlier are common to every stochastic rounding scheme. Since all such schemes give offspring counts with the same marginal distributions, properties such as the marginal offspring variance are common to all stochastic roundings. Indeed it is easy to see that the marginal variance of the offspring counts, $\text{Var}[\nu_t^{(i)} \mid w_t^{(i)}]$ is as small as possible under the constraint of unbiasedness, and as such this is sometimes referred to as minimal-variance resampling. By definition, the support of an offspring count $\nu_t^{(i)}$ given that the associated weight lies in the interval $K/N \leq w_t^{(i)} < (K+1)/N$ is $\{K, K+1\}$. All stochastic roundings are also degenerate when the weights are all equal, i.e. $w_t^{(1:N)} = (1, \dots, 1)/N$ implies $\nu_t^{(1:N)} = (1, \dots, 1)$ almost surely.

2.5 Conditional SMC ✓

⟨sec:condSMC⟩

Be consistent with upper/lower case letters: X, x etc.

Andrieu, Doucet, and Holenstein (2010) propose a number of “particle MCMC” algorithms, which combine SMC with MCMC in order to improve performance in certain situations. One of their algorithms, the *particle Gibbs* sampler (Andrieu, Doucet, and Holenstein 2010, Section 2.4.3), is of particular interest in the current work. For one thing, genealogies are particularly critical to its performance, and for another, the particle step uses a variant SMC algorithm which alters the distribution of genealogies.

In this section, we first introduce the particle Gibbs algorithm and the conditional SMC update, then discuss how ancestral degeneracy impacts the performance of particle Gibbs and how ancestor sampling mitigates this.

2.5.1 Particle Gibbs

To motivate the particle Gibbs algorithm, we introduce a parametrised state space model and explain how combining SMC updates with MCMC sampling allows us to tackle the related inferences effectively. The particle Gibbs algorithm can be applied much more broadly, but this application is particularly intuitive and exhibits all the features of interest to our genealogical study.

Consider a parametrised state space model of the form

$$\begin{aligned}\theta &\sim p(\cdot) \\ X_0 &\sim \mu^\theta(\cdot) \\ X_{t+1} \mid X_t &\sim K_{t+1}^\theta(\cdot \mid X_t) && \text{for } t = 0, \dots, T-1 \\ Y_t \mid X_t &\sim g_t^\theta(\cdot \mid X_t) && \text{for } t = 0, \dots, T\end{aligned}$$

exactly like (2.1) except that the specification is now parametrised by θ (which may be multi-dimensional), and we place a prior distribution on θ . As usual, p , μ^θ , (K_t^θ) and (g_t^θ) are part of the model and are assumed to be known but not necessarily tractable.

Suppose that, given some data $y_{0:T}$, we wish to generate Monte Carlo samples from the joint posterior distribution of $X_{0:T}$ and θ . (Even if we are only interested in inferring θ , for instance, it is often more practical to target the joint posterior and then marginalise.) Notice that we are now working with a finite time horizon $T \in \mathbb{N}$. The inference of interest here is not inherently sequential; we are building an MCMC algorithm to sample from a single target distribution which happens to include some sequentially correlated components.

The conditional dependence structure of the model invites the use of a (partially-collapsed) Gibbs sampler, sampling alternately from the conditional distributions $p(\theta \mid x_{0:T}, y_{0:T})$ and $p(x_{0:T} \mid \theta, y_{0:T})$. The θ update,

$$p(d\theta \mid x_{0:T}, y_{0:T}) \propto p(d\theta)p(x_{0:T}, y_{0:T} \mid \theta),$$

is often quite straightforward, if not analytically then by employing a Metropolis-Hastings step based on the current sampled values of θ and $x_{0:T}$. The X update, meanwhile, is high-dimensional with strong sequential correlations: exactly the situation in which one might use SMC. For the X update, we need a sample from

$$p(dx_{0:T} \mid \theta, y_{0:T}) \propto \mu^\theta(dx_0)g_0^\theta(y_0 \mid x_0) \prod_{s=1}^T K_s^\theta(dx_s \mid x_{s-1})g_s^\theta(y_s \mid x_s), \quad (2.16) \quad \boxed{\text{eq:PG_Xpost}}$$

which can be obtained by running an SMC smoother then sampling one trajectory from its output in proportion to the associated weight.

However, the Markov chain associated to the procedure just described does not admit $p(x_{0:T}, \theta \mid y_{0:T})$ as an invariant distribution. It *approximately* targets this distribution,

2 Background

with some bias. A Gibbs sampler targeting $p(x_{0:T}, \theta \mid y_{0:T})$ exactly can be constructed by replacing the SMC step with a *conditional SMC* step, which takes into account the value of $x_{0:T}$ sampled at the previous iteration, as well as the observations and the current value of θ .

A conditional SMC algorithm for this scenario is presented in Algorithm 2. In contrast to Algorithm 1, the input now includes $x_{0:T}^*$ and $a_{0:T}^*$, which encode the states and parental indices, respectively, of the *immortal trajectory* (so called because it “survives” the SMC run with probability one). Within a particle Gibbs algorithm, the immortal trajectory is set to the trajectory sampled at the previous iteration. The resampling step now assigns the immortal offspring to the immortal parent deterministically, and the state of the immortal particle is also updated deterministically rather than via the Markov kernel. As in standard SMC, there is a choice of RESAMPLE procedures, but some care is needed to ensure the correct treatment of the immortal particle (Lee, Murray, and Johansen 2019).

Input: $T, N, \mu^\theta, (K_t^\theta), (g_t^\theta), y_{0:T}, x_{0:T}^*, a_{0:T}^*$
Set $X_0^{(a_0^*)} \leftarrow x_0^*$
for $i \in \{1, \dots, N\} \setminus a_0^*$ **do** Sample $X_0^{(i)} \sim \mu(\cdot)$
for $i \in \{1, \dots, N\}$ **do** $w_0^{(i)} \leftarrow \left\{ \sum_{j=1}^N g_0^\theta(y_0 \mid X_0^{(j)}) \right\}^{-1} g_0^\theta(y_0 \mid X_0^{(i)})$
for $t \in \{1, \dots, T\}$ **do**
 Set $a_{t-1}^{(a_t^*)} \leftarrow a_{t-1}^*, X_t^{(a_t^*)} \leftarrow x_t^*$
 Sample $a_{t-1}^{(1:N)} \setminus a_{t-1}^* \sim \text{RESAMPLE}(\{1, \dots, N\}, w_{t-1}^{(1:N)})$
 for $i \in \{1, \dots, N\} \setminus a_t^*$ **do** Sample $X_t^{(i)} \sim K_t^\theta(\cdot \mid X_{t-1}^{(a_{t-1}^{(i)})})$
 for $i \in \{1, \dots, N\}$ **do** $w_t^{(i)} \leftarrow \left\{ \sum_{j=1}^N g_t^\theta(y_t \mid X_t^{(j)}) \right\}^{-1} g_t^\theta(y_t \mid X_t^{(i)})$
end

(alg:condSMC) **Algorithm 2:** Conditional sequential Monte Carlo for a parametrised state space model. The immortal particle at each generation has its new state and parental index set deterministically according to the values of $x_{0:T}^*$ and $a_{0:T}^*$ given as input.

The complete particle Gibbs algorithm for this example then consists of alternately sampling from the full conditional distribution of θ (e.g. using a Metropolis-Hastings update) and sampling a trajectory $(x_{0:T}, a_{0:T})$ using conditional SMC. See Andrieu, Doucet, and Holenstein (2010, Section 2.4.3) for more details.

2.5.2 Ancestral degeneracy in particle Gibbs

We have seen in Section 2.3 that the phenomenon of ancestral degeneracy can severely affect the performance of SMC algorithms, particularly in smoothing applications. The SMC update of particle Gibbs is a smoothing problem, however it requires only one sampled trajectory from the smoothing distribution, so one might imagine that we are safe from the curse of ancestral degeneracy. In fact, the loss of ancestors causes a different problem

2 Background

for particle Gibbs: it prevents some components of the Markov chain being refreshed, so that the chain mixes slowly.

To see this, consider the illustration in Figure 2.10, which shows the smoothing trajectories generated by a conditional SMC update at some iteration r . The thick black line is the immortal trajectory given as input, that is, the trajectory sampled by the conditional SMC update at iteration $r - 1$. Backwards in time, the sampled trajectories quickly coalesce until at time 20 all of the trajectories have coalesced. The common trajectory from time 0 to 20 must necessarily be part of the immortal trajectory. A new trajectory (highlighted in purple) is then sampled among the N generated trajectories. Whichever trajectory we sample, it will certainly overlap with the previously sampled trajectory at least from time 0 to 20.

At the next iteration the newly sampled trajectories will again coalesce onto the immortal trajectory, and this behaviour is repeated. If T is too large with respect to N , the early part of the trajectory will very rarely be updated, so the corresponding states will mix very slowly. For further intuition on this phenomenon the reader is directed to Lindsten and Schön (2013, Section 5.4) which provides a very clear exposition.

The meaning of T “too large” here depends on the model and the type of SMC update used, but typically T is determined by the application and N is limited by computational resources, so we may not be able to control their relative size. The other brute-force approach would be to increase the number of iterations of the MCMC algorithm, but this too is infeasible on a limited computational budget. It is therefore worth investing some effort to find an alternative solution to the problem of ancestral degeneracy within particle Gibbs.

2.5.3 Ancestor sampling

(sec:ancsamp) An effective solution (where it is possible to implement it) was proposed by Whiteley (2010) and is known as *ancestor sampling*. It consists of a simple modification to the resampling step within the conditional SMC algorithm. In the basic algorithm with multinomial resampling, at each time step the non-immortal particles are resampled by multinomial resampling according to their weights, while the immortal offspring is deterministically assigned to the immortal parent. That is, at time t , for each $i \in \{1, \dots, N\}$,

$$\mathbb{P}[a_t^{(j)} = i \mid X_{0:t}^{(1:N)}, x_{0:T}^*, a_{0:T}^*] \propto \begin{cases} w_t^{(i)} & j \text{ non-immortal} \\ \mathbb{1}_{\{i=a_t^*\}} & j \text{ immortal.} \end{cases}$$

Ancestor sampling combines the resampling step with a backward simulation step for the immortal particle. Instead of deterministically inheriting the “correct” parent, the immortal particle samples its parent among all N possible parents. This is justified in the same way as backwards simulation in general (Section ??), provided the ancestor sampling probabilities are chosen correctly, although we now apply the backwards sampling step to the immortal trajectory only. Ancestor sampling can also be implemented for other

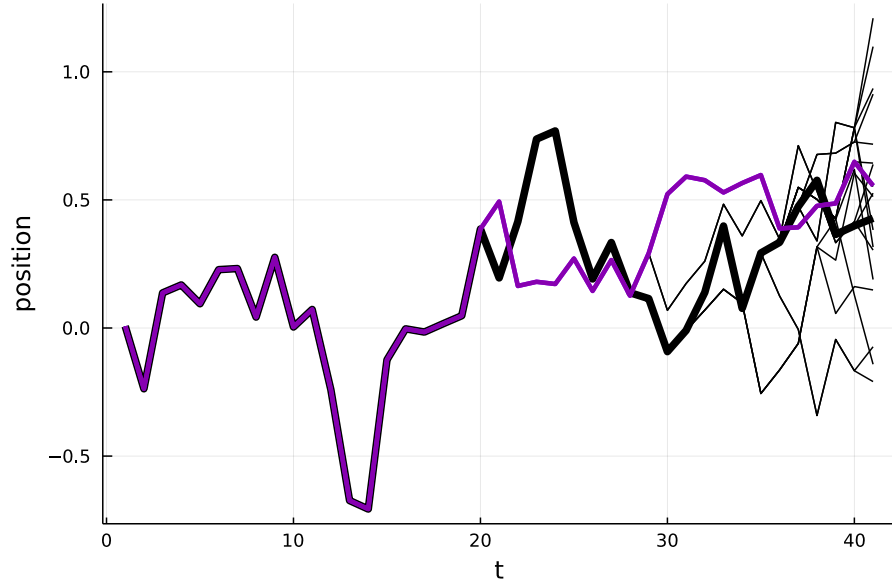


Figure 2.10: Illustration of how ancestral degeneracy causes particle Gibbs to mix slowly on some components. The thick black line is the immortal trajectory, i.e. the sampled trajectory from the previous iteration. Other lines are all of the trajectories generated by conditional SMC. One of these (highlighted in purple) is the sampled trajectory at the current iteration. Due to ancestral degeneracy, the current sample (purple) coincides with the previous sample (thick black) up to time 20, so the components $x_{0:20}$ are not updated in this iteration.

`<fig:PG_ancdegen>`

2 Background

choices of RESAMPLE, using the same backward simulation probabilities (but of course the resampling probabilities for non-immortal particles will be different, and there may be some additional dependence between parental indices). For simplicity we here restrict ourselves to multinomial resampling.

Input: $T, N, \mu^\theta, (K_t^\theta), (g_t^\theta), y_{0:T}, x_{0:T}^*, a_{0:T}^*$
Set $X_0^{(a_0^*)} \leftarrow x_0^*$
for $i \in \{1, \dots, N\} \setminus a_0^*$ **do** Sample $X_0^{(i)} \sim \mu^\theta(\cdot)$
for $i \in \{1, \dots, N\}$ **do** $w_0^{(i)} \leftarrow \left\{ \sum_{j=1}^N g_0^\theta(y_0 \mid X_0^{(j)}) \right\}^{-1} g_0^\theta(y_0 \mid X_0^{(i)})$
for $t \in \{1, \dots, T\}$ **do**
 Set $X_t^{(a_t^*)} \leftarrow x_t^*$
 Sample $a_{t-1}^{(a_t^*)} \sim \text{Categorical}(\{1, \dots, N\}, w_{t-1}^{(1:N)} q_t^\theta(x_t^* \mid X_{t-1}^{(1:N)}))$
 Sample $a_{t-1}^{(1:N)} \setminus a_{t-1}^{(a_t^*)} \sim \text{RESAMPLE}(\{1, \dots, N\}, w_{t-1}^{(1:N)})$
 for $i \in \{1, \dots, N\} \setminus a_t^*$ **do** Sample $X_t^{(i)} \sim q_t^\theta(\cdot \mid X_{t-1}^{(a_{t-1}^{(i)})})$
 for $i \in \{1, \dots, N\}$ **do** $w_t^{(i)} \leftarrow \left\{ \sum_{j=1}^N g_t^\theta(y_t \mid X_t^{(j)}) \right\}^{-1} g_t^\theta(y_t \mid X_t^{(i)})$
end

`<alg:condSMC_ancsamp>` **Algorithm 3:** Conditional sequential Monte Carlo with ancestor sampling for a parametrised state space model. The parent of the “immortal particle” is updated at each iteration via an on-line backward simulation step. The second parameter of the Categorical variable should be interpreted element-wise.

Assume that the smoothing distributions admit densities, that is $\mu^\theta(\cdot)$ and $K_t^\theta(\cdot \mid x)$ admit densities for all x, t . Denote the density of K_t^θ by q_t^θ . Define the trajectories $X_{t,0:t}^{(i)}$ (for any t, i) as in Section 2.1.4, starting from $X_{t,t}^{(i)} := X_t^{(i)}$ and tracing back the states of the parents via $X_{t,s}(i) = X_{t,s+1}^{(a_{t,s}^{(i)})}$. Then the correct resampling probabilities are, for each i ,

$$\mathbb{P}[a_t^{(j)} = i \mid X_{0:t}^{(1:N)}, x_{0:T}^*, a_{0:T}^*] \propto \begin{cases} w_t^{(i)} & j \text{ non-immortal} \\ w_t^{(i)} \frac{p((X_{t,0:t}^{(i)}, x_{t+1:T}^*) \mid \theta, y_{0:T})}{p(X_{t,0:t}^{(i)} \mid \theta, y_{0:t})} & j \text{ immortal.} \end{cases} \quad (2.17) \quad \boxed{\text{eq:ancsamp_1}}$$

The ratio of densities can be interpreted as the conditional probability that the whole trajectory is the concatenation of $X_{t,0:t}^{(i)}$ with $x_{t+1:T}^*$, given that its first $t+1$ states are $X_{t,0:t}^{(i)}$. To simplify the ratio, use (2.18) to write

$$p(X_{t,0:t}^{(i)} \mid \theta, y_{0:t}) \propto \mu^\theta(X_{t,0}^{(i)}) g_0^\theta(y_0 \mid X_{t,0}^{(i)}) \prod_{s=1}^t q_s^\theta(X_{t,s}^{(i)} \mid X_{t,s-1}^{(i)}) g_s^\theta(y_s \mid X_{t,s}^{(i)})$$

2 Background

and

$$p((X_{t,0:t}^{(i)}, x_{t+1:T}^*) \mid \theta, y_{0:T}) \propto \mu^\theta(X_{t,0}^{(i)}) g_0^\theta(y_0 \mid X_{t,0}^{(i)}) \left\{ \prod_{s=1}^t q_s^\theta(X_{t,s}^{(i)} \mid X_{t,s-1}^{(i)}) g_s^\theta(y_s \mid X_{t,s}^{(i)}) \right\} \\ \times q_{t+1}^\theta(x_{t+1}^* \mid X_{t,t}^{(i)}) g_{t+1}^\theta(y_{t+1} \mid x_{t+1}^*) \left\{ \prod_{s=t+2}^T q_s^\theta(x_s^* \mid x_{s-1}^*) g_s^\theta(y_s \mid x_s^*) \right\}.$$

The ratio then becomes

$$\frac{p((X_{t,0:t}^{(i)}, x_{t+1:T}^*) \mid \theta, y_{0:T})}{p(X_{t,0:t}^{(i)} \mid \theta, y_{0:t})} \propto q_{t+1}^\theta(x_{t+1}^* \mid X_{t,t}^{(i)}) g_{t+1}^\theta(y_{t+1} \mid x_{t+1}^*) \prod_{s=t+2}^T q_s^\theta(x_s^* \mid x_{s-1}^*) g_s^\theta(y_s \mid x_s^*) \\ \propto q_{t+1}^\theta(x_{t+1}^* \mid X_{t,t}^{(i)}) \\ = q_{t+1}^\theta(x_{t+1}^* \mid X_t^{(i)}).$$

The probabilities in (2.19) become

$$\mathbb{P}[a_t^{(j)} = i \mid X_{0:t}^{(1:N)}, x_{0:T}^*, a_{0:T}^*] \propto \begin{cases} w_t^{(i)} & j \text{ non-immortal} \\ w_t^{(i)} q_{t+1}^\theta(x_{t+1}^* \mid X_t^{(i)}) & j \text{ immortal.} \end{cases} \quad (2.18) \quad \boxed{\text{eq:ancsamp-}}$$

The conditional SMC algorithm with this adaptation is presented in Algorithm 3.

We see that, in order to do ancestor sampling, we need a stronger assumption on the Markov kernels than was required to simply run the conditional SMC algorithm: we now require that, for each t , K_t^θ admits a density q_t^θ and that $q_t^\theta(\cdot \mid x)$ can be evaluated pointwise for any x , whereas previously we only needed to draw samples from $K_t^\theta(\cdot \mid x)$ for any x . This additional requirement rules out ancestor sampling in some applications.

Recall that the usual backward simulation procedure requires a full forward pass to calculate the future states before the backward simulation probabilities can be computed. Ancestor sampling, on the other hand, does not require a forward pass because it only computes backward simulation probabilities for the immortal trajectory, for which all the future states are known in advance. This means that the additional computational cost of implementing ancestor sampling is negligible.

Why ancestor sampling works

We know that complete backward simulation eradicates ancestral degeneracy by breaking up the lineages into disjoint pieces (Section ??). But here we are only backward-simulating one of the N particles, leaving the other $N - 1$ lineages to coalesce as usual. So how does this help?

Recall that in particle Gibbs ancestral degeneracy is not itself a problem, because we only require a single sample from the smoothing distribution. The problem is that the consecutive samples are highly correlated, because of the repeated coalescence onto the immortal lineage. The contribution of ancestor sampling is to break up the immortal

2 Background

trajectory so that it no longer appears in the sampled lineages; see Figure 2.11. While the non-immortal trajectories may still coalesce, they no longer preferentially coalesce onto the immortal trajectory. In turn, the sampled trajectory that is output does not overlap unduly with the immortal trajectory that was the previous output, and this completely solves the problem of the slow-mixing particle Gibbs chain.

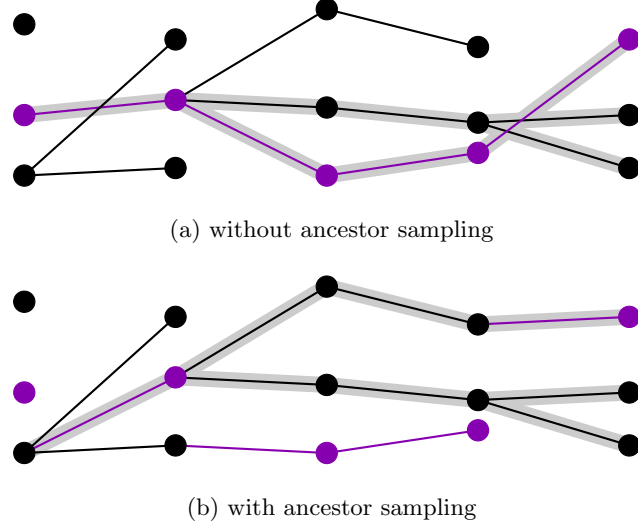


Figure 2.11: Illustration of how ancestor sampling prevents coalescence onto the immortal trajectory. Immortal particles are highlighted in purple, along with their parent-offspring edges (the given ones in (a) and the ancestor-sampled ones in (b)). The resulting lineages of the terminal particles are highlighted in grey. In (a), the lineages of the terminal particles coalesce onto the immortal trajectory. Imagine time stretching further back: the lineages would continue to coincide with the immortal trajectory forever. In (b), the lineages still coalesce, but not onto the immortal trajectory. The immortal trajectory no longer exists as a lineage.

`<fig:whyASworks>`

Bibliography

- Andrieu, Christophe, Arnaud Doucet, and Roman Holenstein (2010). “Particle Markov Chain Monte Carlo Methods”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 72.3, pp. 269–342.
- Baum, Leonard E. et al. (1970). “A Maximization Technique Occurring in the Statistical Analysis of Probabilistic Functions of Markov Chains”. In: *The Annals of Mathematical Statistics* 41, pp. 164–171.
- Berestycki, Nathanaël (2009). *Recent Progress in Coalescent Theory*. 0909.3985. ArXiv.
- Cannings, Chris (1974). “The Latent Roots of Certain Markov Chains Arising in Genetics: A New Approach, I. Haploid Models”. In: *Advances in Applied Probability* 6.2, pp. 260–290.
- (1975). “The Latent Roots of Certain Markov Chains Arising in Genetics: A New Approach, II. Further Haploid Models”. In: *Advances in Applied Probability* 7.2, pp. 264–282.
- Carpenter, James, Peter Clifford, and Paul Fearnhead (1999). “Improved Particle Filter for Nonlinear Problems”. In: *IEE Proceedings — Radar, Sonar and Navigation* 146.1, pp. 2–7.
- Chopin, Nicolas and Omiros Papaspiliopoulos (2020). *An Introduction to Sequential Monte Carlo*. Springer.
- Corenflos, Adrien et al. (2021). *Differentiable Particle Filtering via Entropy-Regularized Optimal Transport*. Tech. rep. 2102.07850. ArXiv.
- Crisan, Dan, Pierre Del Moral, and Terry Lyons (1999). “Discrete Filtering Using Branching and Interacting Particle Systems”. In: *Markov Processes and Related Fields* 5.3, pp. 293–318.
- Crisan, Dan and Terry Lyons (1997). “Nonlinear Filtering and Measure-Valued Processes”. In: *Probability Theory and Related Fields* 109.2, pp. 217–244.
- (1999). “A Particle Approximation of the Solution of the Kushner–Stratonovitch Equation”. In: *Probability Theory and Related Fields* 115.4, pp. 549–578.
- Del Moral, Pierre (2004). *Feynman–Kac Formulae: Genealogical and Interacting Particle Systems with Applications*. Springer.
- (2013). *Mean Field Simulation for Monte Carlo Integration*. Chapman and Hall/CRC.
- Douc, Randal, Olivier Cappé, and Eric Moulines (2005). “Comparison of Resampling Schemes for Particle Filtering”. In: *Proceedings of the 4th International Symposium on Image and Signal Processing and Analysis*. IEEE, pp. 64–69.

Bibliography

- Doucet, Arnaud and Adam M. Johansen (2011). “A Tutorial on Particle Filtering and Smoothing: Fifteen Years Later”. In: *Handbook of Nonlinear Filtering*. OUP, pp. 656–704.
- Durrett, Richard (2008). *Probability Models for DNA Sequence Evolution*. Springer Science & Business Media.
- Efron, Bradley and Robert J. Tibshirani (1994). *An Introduction to the Bootstrap*. CRC press.
- Etheridge, Alison (2011). *Some Mathematical Models from Population Genetics: École D’Été de Probabilités de Saint-Flour XXXIX-2009*. Springer.
- Evensen, Geir (1994). “Sequential Data Assimilation with a Nonlinear Quasi-Geostrophic Model Using Monte Carlo Methods to Forecast Error Statistics”. In: *Journal of Geophysical Research: Oceans* 99.C5, pp. 10143–10162.
- Fearnhead, Paul and Peter Clifford (2003). “On-line Inference for Hidden Markov Models via Particle Filters”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 65.4, pp. 887–899.
- Fisher, Ronald Aylmer (1923). “On the Dominance Ratio”. In: *Proceedings of the Royal Society of Edinburgh* 42, pp. 321–341.
- (1930). “The Distribution of Gene Ratios for Rare Mutations”. In: *Proceedings of the Royal Society of Edinburgh* 50, pp. 205–220.
- Gerber, Mathieu, Nicolas Chopin, and Nick Whiteley (2019). “Negative Association, Ordering and Convergence of Resampling Methods”. In: *The Annals of Statistics* 47.4, pp. 2236–2260.
- Gordon, Neil J., David J. Salmond, and Adrian F. M. Smith (1993). “Novel Approach to Nonlinear/Non-Gaussian Bayesian State Estimation”. In: *IEEE Proceedings F (Radar and Signal Processing)*. Vol. 140. 2. IET, pp. 107–113.
- Hol, Jeroen D. (2004). “Resampling in Particle Filters”. LiTH-ISY-EX-ET-0283-2004. Internship report. Linköping University.
- Hol, Jeroen D., Thomas B. Schön, and Fredrik Gustafsson (2006). “On Resampling Algorithms for Particle Filters”. In: *Nonlinear Statistical Signal Processing Workshop*. IEEE, pp. 79–82.
- Huang, Chaofan, Vengazhiyil Roshan Joseph, and Simon Mak (2020). *Population Quasi-Monte Carlo*. Tech. rep. 2012.13769. ArXiv.
- Jacob, Pierre E., Lawrence M. Murray, and Sylvain Rubenthaler (2015). “Path Storage in the Particle Filter”. In: *Statistics and Computing* 25.2, pp. 487–496.
- Jazwinski, Andrew H. (2007). *Stochastic Processes and Filtering Theory*. Courier Corporation.
- Joag-Dev, Kumar and Frank Proschan (1983). “Negative Association of Random Variables with Applications”. In: *The Annals of Statistics*, pp. 286–295.
- Kalman, Rudolph Emil (1960). “A New Approach to Linear Filtering and Prediction Problems”. In: *Journal of Basic Engineering* 82.1, pp. 35–45.

Bibliography

- Kingman, John F. C. (1982a). “Exchangeability and the Evolution of Large Populations”. In: *Proceedings of the International Conference on Exchangeability in Probability and Statistics, Rome, 6th-9th April, 1981, in Honour of Professor Bruno de Finetti*. North-Holland, Amsterdam.
- (1982b). “On the Genealogy of Large Populations”. In: *Journal of Applied Probability* 19.A, pp. 27–43.
- (1982c). “The Coalescent”. In: *Stochastic Processes and Their Applications* 13.3, pp. 235–248.
- Kitagawa, Genshiro (1996). “Monte Carlo Filter and Smoother for Non-Gaussian Non-linear State Space Models”. In: *Journal of Computational and Graphical Statistics* 5.1, pp. 1–25.
- Koskela, Jere et al. (2018). *Asymptotic genealogies of interacting particle systems with an application to sequential Monte Carlo*. Mathematics e-print 1804.01811. ArXiv.
- Kuipers, Lauwerens and Harald Niederreiter (1974). *Uniform Distribution of Sequences*. John Wiley & Sons.
- Lee, Anthony, Lawrence M. Murray, and Adam M. Johansen (2019). “Resampling in Conditional SMC Algorithms”. Unpublished.
- Lindsten, Fredrik and Thomas B. Schön (2013). “Backward Simulation Methods for Monte Carlo Statistical Inference”. In: *Foundations and Trends in Machine Learning* 6.1, pp. 1–143.
- Liu, Jun S. and Rong Chen (1995). “Blind Deconvolution via Sequential Imputations”. In: *Journal of the American Statistical Association* 90.430, pp. 567–576.
- (1998). “Sequential Monte Carlo Methods for Dynamic Systems”. In: *Journal of the American Statistical Association* 93.443, pp. 1032–1044.
- Liu, Jun S., Rong Chen, and Tanya Logvinenko (2001). “A Theoretical Framework for Sequential Importance Sampling with Resampling”. In: *Sequential Monte Carlo Methods in Practice*. Springer, pp. 225–246.
- Moran, Patrick Alfred Pierce (1958). “Random Processes in Genetics”. In: *Mathematical Proceedings of the Cambridge Philosophical Society*. Vol. 54. 1. Cambridge University Press, pp. 60–71.
- Myers, Aaron et al. (2021). “Sequential Ensemble Transform for Bayesian Inverse Problems”. In: *Journal of Computational Physics* 427, p. 110055.
- Pitt, Michael K. and Neil Shephard (1999). “Filtering via Simulation: Auxiliary Particle Filters”. In: *Journal of the American Statistical Association* 94.446, pp. 590–599.
- Rauch, Herbert E., C. T. Striebel, and F. Tung (1965). “Maximum Likelihood Estimates of Linear Dynamic Systems”. In: *AIAA Journal* 3.8, pp. 1445–1450.
- Reich, Sebastian (2013). “A Nonparametric Ensemble Transform Method for Bayesian Inference”. In: *SIAM Journal on Scientific Computing* 35.4, A2013–A2024.
- Saunders, Ian W., Simon Tavaré, and G. A. Watterson (1984). “On the Genealogy of Nested Subsamples from a Haploid Population”. In: *Advances in Applied Probability*, pp. 471–491.

Bibliography

- Spouge, John L. (2014). “Within a Sample from a Population, the Distribution of the Number of Descendants of a Subsample’s Most Recent Common Ancestor”. In: *Theoretical Population Biology* 92, pp. 51–54.
- Srinivasan, Aravind (2001). “Distributions on Level-Sets with Applications to Approximation Algorithms”. In: *Proceedings of the 42nd IEEE Symposium on Foundations of Computer Science*. IEEE, pp. 588–597.
- Vidoni, Paolo (1999). “Exponential Family State Space Models Based on a Conjugate Latent Process”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 61.1, pp. 213–221.
- Wakeley, John (2009). *Coalescent Theory: An Introduction*. Roberts & Co. Publishers.
- Wan, Eric A. and Rudolph van der Merwe (2000). “The Unscented Kalman Filter for Nonlinear Estimation”. In: *Proceedings of the IEEE 2000 Adaptive Systems for Signal Processing, Communications, and Control Symposium*. IEEE, pp. 153–158.
- Whiteley, Nick (2010). “Discussion on ‘Particle Markov Chain Monte Carlo Methods’”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 72.3, pp. 306–307.
- Whitley, Darrell (1994). “A Genetic Algorithm Tutorial”. In: *Statistics and Computing* 4.2, pp. 65–85.
- Wright, Sewall (1931). “Evolution in Mendelian Populations”. In: *Genetics* 16.2, pp. 97–159.