

A Tutorial on Particle Filters

Maarten Speekenbrink*

*Experimental Psychology
University College London*

Abstract

This tutorial aims to provide an accessible introduction to particle filters, and sequential Monte Carlo (SMC) more generally. These techniques allow for Bayesian inference in complex dynamic state-space models and have become increasingly popular over the last decades. The basic building blocks of SMC – sequential importance sampling and resampling – are discussed in detail with illustrative examples. A final example presents a particle filter for estimating time-varying learning rates in a probabilistic category learning task.

Keywords: Particle filter, Sequential Monte Carlo, State-space model, Sequential Bayesian inference

1 Particle filters, and sequential Monte Carlo (SMC) techniques more gen-
2 erally, are a class of simulation-based techniques which have become in-
3 creasingly popular over the last decades to perform Bayesian inference in
4 complex dynamic statistical models (e.g., Doucet, de Freitas, and Gordon,
5 2001b; Doucet and Johansen, 2011). Particle filters are generally applied

*Department of Experimental Psychology, University College London, Gower Street, London WC1E 6BT, England.

Email address: m.speekenbrink@ucl.ac.uk (Maarten Speekenbrink)

URL: www.ucl.ac.uk/speekenbrink-lab/ (Maarten Speekenbrink)

6 to so-called filtering problems, where the objective is to estimate the latent
7 states of a stochastic process on-line, such that, after each sequential obser-
8 vation, the state giving rise to that observation is estimated. For instance,
9 in a category learning task, we might want to infer how people use the
10 features of objects to categorize them. Due to learning, we would expect
11 their categorization strategy to change over time. Traditionally, a formal
12 learning model such as ALCOVE (Kruschke, 1992) would be used for this
13 purpose, which describes how feedback on their categorization decisions af-
14 fects people’s momentary strategy. However, these models usually assume
15 a deterministic updating process, which may be too restrictive. Ideally, we
16 would like to estimate someone’s strategy – which we can view as the la-
17 tent state of their decision process – from trial to trial whilst allowing for
18 stochastic transitions between states. Estimating the current categoriza-
19 tion strategy is a difficult task, however, as a single categorization decision
20 at each point in time provides relatively little information about people’s
21 complete categorization strategy, i.e. their potential categorizations of all
22 possible stimuli. Assuming trial-to-trial changes to a state (strategy) are
23 noisy but relatively small, we may however be able to gain some insight
24 into the current state from all previous categorization decisions someone
25 made. This filtering problem is generally not analytically tractable; ana-
26 lytical results are only available for the restricted class of linear Gaussian
27 state-space models. As particle filters are applicable to the much broader
28 class of non-linear non-Gaussian state-space models, they open up interest-
29 ing possibilities to study a broad range of dynamic processes in psychology.

30 A graphical representation of a generic particle filter (see Section 4.3) is
31 given in Figure 1. Particle filters operate on a set of randomly sampled values
32 of a latent state or unknown parameter. The sampled values, generally

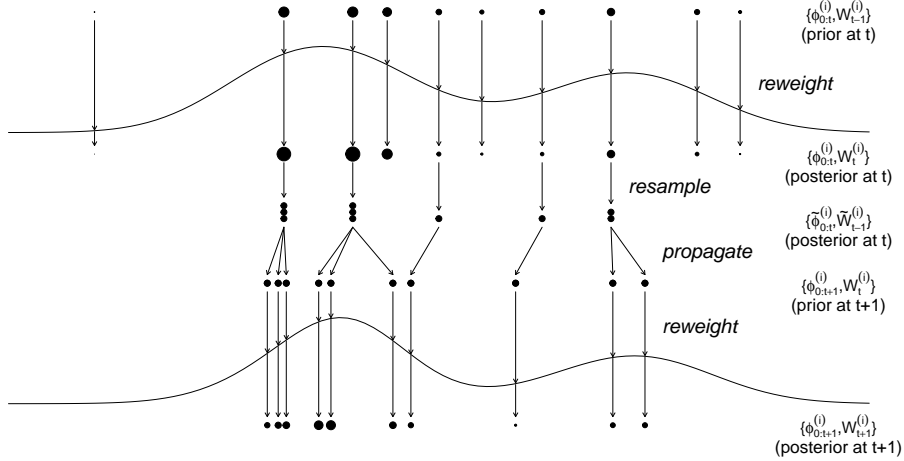


Figure 1: Schematic representation of a generic particle filter (after Doucet et al., 2001a). Standing at time t , we have a set of weighted particles $\{\phi_{0:t}^{(i)}, W_{t-1}^{(i)}\}$ representing the prior distribution at t . Each particle $\phi_{0:t}^{(i)}$ is a multidimensional variable which represents the whole path of the latent state from time 0 up to the current time point t , such that each dimension represents the value of the state at a particular time point. The location of the dots in the graph reflect $\phi_t^{(i)}$, the value of the state at the current time point, i.e. the dimension of each particle reflecting the current state. The size of each dot reflects the weight $W_{t-1}^{(i)}$ (“prior at t ”). In the *reweight* step, the weights are updated to $W_t^{(i)}$ partly as a function of $p(y_t|\phi_t^{(i)})$, the likelihood of observation y_t according to each sampled state value $\phi_t^{(i)}$ (solid line). The resulting set $\{\phi_{0:t}^{(i)}, W_t^{(i)}\}$ of weighted particles approximates the posterior distribution (“posterior at t ”) of the latent state paths. The *resampling* step duplicates values $\phi_{0:t}^{(i)}$ with high weights $W_t^{(i)}$, and eliminates those with low weights, resulting in the set of uniformly weighted particles $\{\tilde{\phi}_{0:t}^{(i)}, \tilde{W}_t^{(i)} = 1/N\}$ which is approximately distributed according to the posterior (second “posterior at t ”). In the *propagate* step, values of states $\phi_{t+1}^{(i)}$ at the next time point are sampled and added to each particle to account for state transitions, forming a prior distribution for time $t+1$ (“prior at $t+1$ ”). Thus, at each new time point, the particles grow in dimension because the whole path of the latent state now incorporates the new time point as well. The particles are then reweighted in response to the likelihood of the new observation y_{t+1} to approximate the posterior distribution at $t+1$ (“posterior at $t+1$ ”), etc.

33 referred to as “particles”, are propagated over time to track the posterior
34 distribution of the state or parameter at each point in time. Each particle
35 is assigned a weight in relation to its posterior probability. To increase their
36 accuracy, SMC techniques resample useful particles from the set according to
37 these weights. This resampling introduces interaction between the particles,
38 and the term “interacting particle filters” was coined by Del Moral (1996),
39 who showed how the method relates to techniques used in physics to analyse
40 the movement of particles.

41 Particle filters have successfully solved difficult problems in machine
42 learning, such as allowing robots to simultaneously map their environment
43 and localize their position within it (Montemerlo, Thrun, Koller, and Weg-
44 breit, 2002), and the automated tracking of multiple objects in naturalistic
45 videos (Isard and Blake, 1998; Nummiaro, Koller-Meier, and Gool, 2003).
46 More recently, particle filters have also been proposed as models of human
47 cognition, for instance how people learn to categorize objects (Sanborn, Grif-
48 fiths, and Navarro, 2010), how they detect and predict changes (Brown and
49 Steyvers, 2009) as well as make decisions (Yi, Steyvers, and Lee, 2009) in
50 changing environments.

51 The aim of this tutorial is to provide readers with an accessible introduc-
52 tion to particle filters and SMC. We will discuss the foundations of SMC,
53 sequential importance sampling and resampling, in detail, using simple ex-
54 amples to highlight important aspects of these techniques. We start with a
55 discussion of importance sampling, which is a Monte Carlo integration tech-
56 nique which can be used to efficiently compute expected values of random
57 variables, including expectations regarding the posterior probabilities of la-
58 tent states or parameters. We will then move on to sequential importance
59 sampling, an extension of importance sampling which allows for efficient

60 computation in sequential inference problems. After introducing resampling
61 as a means to overcome some problems in sequential importance sampling,
62 we have all the ingredients to introduce a generic particle filter. After dis-
63 cussing limitations and extensions of SMC, we will conclude with a more
64 complex example involving the estimation of time-varying learning rates in
65 a probabilistic category learning task.

66 **1. Importance sampling**

67 Importance Sampling (IS) is a Monte Carlo integration technique. It
68 can be used to efficiently solve high-dimensional integration problems when
69 analytical solutions are difficult or unobtainable. In statistics, it is often
70 used to approximate expected values of random variables, which is what we
71 will focus on here. If we have a sample of realizations of a random variable
72 Y , we can estimate the expected value by computing a sample average. We
73 do this when we have data from experiments, and it is also the idea behind
74 basic Monte Carlo integration. Importance sampling is based on the same
75 idea, but rather than sampling values from the true distribution of Y , values
76 are sampled from a different distribution, called the importance distribution.
77 Sampling from a different distribution can be useful to focus more directly
78 on the estimation problem at hand, or if it is problematic to sample from
79 the target distribution. To correct for the fact that the samples were drawn
80 from the importance distribution and not the target distribution, weights
81 are assigned to the sampled values which reflect the difference between the
82 importance and target distribution. The final estimate is then a weighted
83 average of the randomly sampled values.

84 Suppose we wish to compute the expected value of an arbitrary function

85 f of a random variable Y which is distributed according to a probability
86 distribution p :

$$87 \quad \mathbb{E}_p[f(Y)] \triangleq \int f(y)p(y) dy.$$

88 This is just the usual definition of an expected value (we use \mathbb{E}_p to denote
89 an expectation of a random variable with distribution p , and the symbol
90 \triangleq to denote ‘is defined as’). The function f depends on what we want to
91 compute. For instance, choosing $f(y) = y$ would result in computing the
92 mean of Y , while choosing $f(y) = (y - \mathbb{E}_p[f(Y)])^2$ would result in computing
93 the variance of Y . It is often not possible to find an analytical solution to
94 the integral above, in which case we have to turn to some form of numerical
95 approximation. A basic Monte Carlo approximation is to draw a number of
96 independent samples from p and then compute a sample average from these
97 random draws:

98 **Algorithm 1.** Basic Monte Carlo integration for an expected value $\mathbb{E}_p[f(Y)]$

99

- 100 1. (Sample) For $i = 1, \dots, N$, sample $y^{(i)} \sim p(y)$.
- 101 2. (Estimate) Compute the sample average to obtain the Monte Carlo
102 estimate E^{MC} of the expected value:

$$103 \quad E^{\text{MC}} \triangleq \frac{1}{N} \sum_{i=1}^N f(y^{(i)}). \quad (1)$$

104 We let $y^{(i)}$ denote the i -th sampled value and for consistency in terminology,
105 we will refer to these sampled values as “particles” from now on. By the
106 law of large numbers, as the number N of particles approaches infinity, this
107 estimate will converge almost surely¹ to the true value (Robert and Casella,

¹Almost sure convergence means that the probability that the estimate is identical to

108 2004). A limitation of this procedure is that we need to be able to sample
 109 particles according to the distribution p , which is not always possible or
 110 efficient. Importance sampling circumvents this limitation, allowing parti-
 111 cles to be drawn from an arbitrary “instrumental distribution” q . These
 112 particles are then weighted to correct for the fact they were drawn from q
 113 and not the target distribution p . Importance sampling relies on the simple
 114 algebraic identity $a = \frac{a}{b} \times b$ to derive the following *importance sampling*
 115 *fundamental identity* (Robert and Casella, 2004):

$$116 \quad \mathbb{E}_p[f(Y)] = \int \frac{p(y)}{q(y)} q(y) f(y) dy = \mathbb{E}_q[w(Y)f(Y)],$$

117 where we define the importance weight as $w(y) = \frac{p(y)}{q(y)}$. Thus, the expected
 118 value of $f(Y)$ under the target distribution p is identical to the expected
 119 value of the product $w(Y)f(Y)$ under the instrumental distribution q . The
 120 instrumental distribution can be chosen for ease of sampling, or to increase
 121 the efficiency of the estimate (as shown in the example below). The only
 122 restriction on q is that, in the range where $f(y) \neq 0$, q should have the same
 123 support as p (i.e. whenever p assigns non-zero probability to a value y , q
 124 should do so also, so $q(y) > 0$ whenever $p(y) > 0$). More compactly, we can
 125 state this requirement as: if $f(y)p(y) \neq 0$, then $q(y) > 0$. An IS estimate of
 126 the expected value of $f(Y)$ under p is thus obtained by generating a sample
 127 from q and computing a weighted average, as in the following algorithm:

128 **Algorithm 2.** Importance sampling for an expected value $\mathbb{E}_p[f(Y)]$

- 129 1. (Sample) For $i = 1, \dots, N$, sample $y^{(i)} \sim q(y)$.

the true value approaches 1 as N approaches infinity, i.e., $p(\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N f(y^{(i)}) = \mathbb{E}_p(f(Y))) = 1$.

- 130 2. (Weight) For $i = 1, \dots, N$, compute the importance weight $w^{(i)} =$
 131 $\frac{p(y^{(i)})}{q(y^{(i)})}$.
 132 3. (Estimate) Compute a weighted average to obtain the IS estimate:

$$133 \quad E^{\text{IS}} \triangleq \frac{1}{N} \sum_{i=1}^N w^{(i)} f(y^{(i)}). \quad (2)$$

134 As for basic Monte Carlo estimation, the law of large numbers assures that
 135 E^{IS} converges to $\mathbb{E}_p[f(Y)]$ as the number of particles approaches infinity
 136 (Robert and Casella, 2004).

137 It should be stressed that IS is a Monte Carlo integration method, which
 138 can be used to approximate expected values of random variables. It is not
 139 a method to directly generate samples according to the target distribution.
 140 However, we can generate samples which are approximately distributed ac-
 141 cording to the target distribution by resampling particles with replacement
 142 from the set of particles, where we sample a particle $y^{(i)}$ with a probability
 143 proportional to the importance weight $w(y^{(i)})$. This importance sampling
 144 resampling algorithm, which will be discussed in more detail later, can pro-
 145 vide an “empirical” approximation to the distribution p (in the sense that
 146 we use a finite random sample drawn from p to approximate p , just like a his-
 147 togram of observations from an experiment approximates the distribution of
 148 possible observations that could be made in that experiment). We can also
 149 use IS to compute any probability within the distribution p . For instance,
 150 if we want to compute the probability that the value of Y is between a and
 151 b , we can use IS with the indicator function $f(y) = \mathbb{I}(a \geq y \geq b)$, where
 152 the indicator function \mathbb{I} equals 1 when its argument is true and 0 otherwise.
 153 We can do this as it is easy to show that the required probability equals the
 154 expected value of this indicator function: $p(a \geq Y \geq b) = \mathbb{E}_p[\mathbb{I}(a \geq y \geq b)]$.
 155 In practice, the estimated probability is then simply the sum of the impor-

156 tance weights of the particles that lie between a and b . This is illustrated
157 in Figure 2. A few remarks are in order. Firstly, while the estimated prob-
158 abilities are unbiased, in practice, we can only estimate the probability if at
159 least one particle falls within the interval. Secondly, given a set of particles,
160 we can vary the bounds of the interval in between the particles and we will
161 obtain the same estimates, because if two regions capture the same subset
162 of particles, the sum of the weights of those particles will also be identical.
163 Finally, to obtain a precise estimate of a probability, it is wise to tailor the
164 importance distribution to sample solely in the required region, as will be
165 shown in the following example.

166 *1.1. Example: computing the tail probability of the Ex-Gaussian distribution*

167 The ex-Gaussian distribution is a popular distribution to model response
168 times (Van Zandt, 2000). The ex-Gaussian distribution is defined as the sum
169 of an exponential and normal (Gaussian) distributed variable, and has three
170 parameters: μ , σ , and τ , which are respectively the mean and standard
171 deviation of the Gaussian variable, and the rate of the exponential variable.
172 See Figure 3 for an example of an ex-Gaussian distribution.

173 Suppose that for a certain person the distribution of completion times
174 for a task approximately follows an ex-Gaussian distribution with param-
175 eters $\mu = 0.4$, $\sigma = 0.1$ and $\tau = 0.5$, and that we want to know on how
176 many trials that person would fail to complete the task within a time limit
177 of 3 seconds. Looking at Figure 3, we can already see that the probability of
178 non-completion is rather small. As the ex-Gaussian distribution is relatively
179 easy to draw samples from, we can use basic Monte Carlo integration (Algo-
180 rithm 1) to approximate this probability. With a sample size of $N = 2000$,
181 this gave an estimate $p(Y \geq 3) \approx 0.0040$. Compared to the true value,

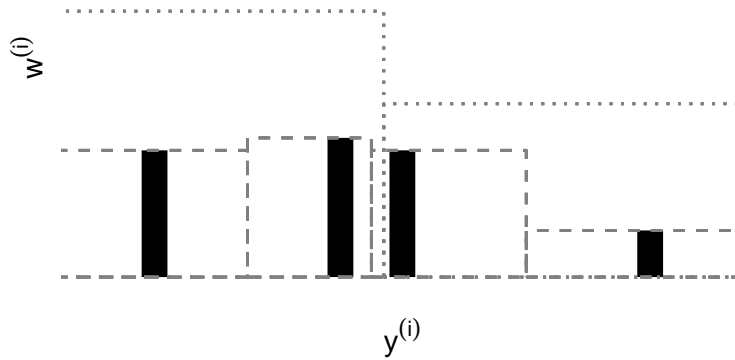


Figure 2: Probability estimation with importance sampling. The locations of the black bars represent particle values ($y^{(i)}$) and the height of the black bars represents the corresponding weights ($w^{(i)}$). The broken and dotted lines represent two different estimates of probabilities from these particles. The broken lines involve smaller regions which each include a single particle, while the dotted lines involve larger regions which include multiple particles. Small changes to the bounds of the regions would leave the estimates unchanged as long as the same particles fall within each region.

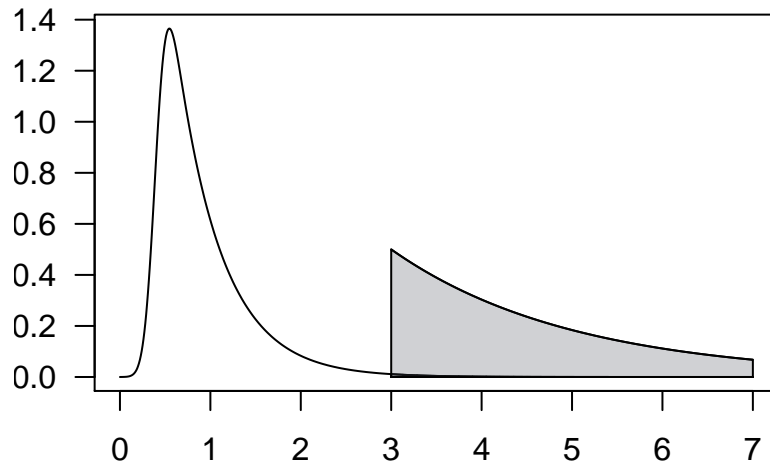


Figure 3: Tail probability estimation for an ex-Gaussian distribution through importance sampling with a shifted Exponential distribution. The solid line with an unshaded region below it reflects the ex-Gaussian distribution. Overlaid and shaded grey is the shifted Exponential distribution which is used as importance distribution.

182 $p(Y \geq 3) = 0.0056$, this estimate is too low by 28.93%. Basic Monte Carlo
183 integration fails here because the exceedance probability $p(Y \geq 3)$ is rela-
184 tively small and we therefore need many samples to obtain an estimate with
185 adequate precision. Because the importance distribution can be tailored to
186 the estimation problem at hand, IS can be much more efficient. To apply
187 IS, we first formulate the desired result as the expected value

$$188 \quad p(Y \geq 3) = \mathbb{E}_p[\mathbb{I}(y \geq 3)].$$

189 We then need to choose an importance distribution. Recall that the only
190 requirement is that the instrumental distribution q has the same support as
191 the target distribution p in the range where $f(y) \neq 0$. For this example, q
192 need thus only be defined over the range $[3; \infty)$. In fact, choosing an impor-
193 tance distribution which does not extend beyond this range is a good idea,
194 because samples outside this range are wasteful as they will not affect the
195 estimate. A reasonable choice is a shifted exponential distribution, shifted
196 to the right to start at 3 rather than 0. With a sample of $N = 2000$ and
197 matching $\tau = 0.5$ to the same value as in the ex-Gaussian distribution, this
198 gives the estimate $p(Y \geq 3) \approx 0.0055$, which deviates from the true value
199 by only 2.63%. This is a representative example and shows the IS estima-
200 tor is much better than the basic Monte Carlo estimator. While using an
201 importance distribution defined over the range $[3, \infty)$, such as the shifted
202 exponential, is a good way to increase the precision of the estimate, we must
203 be careful when choosing the importance distribution. For example, a Nor-
204 mal distribution truncated below at 3 with parameters $\mu = 3$ and $\sigma = 0.1$,
205 resulted in the estimate $P(Y \geq 3) \approx 0.0036$, which is too low by 35.55%
206 and worse than the basic Monte Carlo estimate. The problem with this
207 truncated Normal is that the parameter $\sigma = 0.1$ is set too low, resulting in a

208 distribution with a right tail which is too light compared to the ex-Gaussian
209 distribution. Recall that the importance weights are given by the ratio $\frac{p(y)}{q(y)}$.
210 If the instrumental distribution has lighter tails than the target distribution,
211 there will be relatively few particles that fall in the tails, but for these rare
212 particles, $p(y)$ may be very large compared to $q(y)$, leading to very large
213 weights. In the extreme case, when one or a few importance weights are
214 very large compared to the other weights, the estimate is effectively deter-
215 mined by only one or a few particles, which is obviously bad. For example,
216 in the most extreme estimate resulting from the truncated Normal, there
217 was one particle with a weight of 88.4, while the next highest weight was
218 0.6. For comparison, in the most extreme estimate from the shifted Expo-
219 nential distribution, the largest and second-largest importance weights were
220 both 0.02. Large variation in importance weights results in an estimator
221 with a high variance. This can be clearly seen in Figure 4, which shows
222 the variation in the estimates of the three estimators when applying them
223 repeatedly for 1000 times. While the distribution of the estimates obtained
224 for IS with a shifted exponential distribution is tightly clustered around the
225 true value, both basic Monte Carlo integration and IS with the truncated
226 Normal provide much more variable estimates. Note that these estimates
227 are still unbiased, in the sense that on average, they are equal to the true
228 value. However, the large variance of the estimates means that in practice,
229 we are often quite far off the true value. The positive skew for the truncated
230 Normal shows that IS with this importance distribution underestimates the
231 probability most of the time, but in the rare cases that a particle falls in
232 the right tail, the high weight assigned to that particle results in a large
233 overestimation of the probability. Note that there is no problem with using
234 a truncated Normal per se: increasing the parameter to $\sigma = 1$ results in a

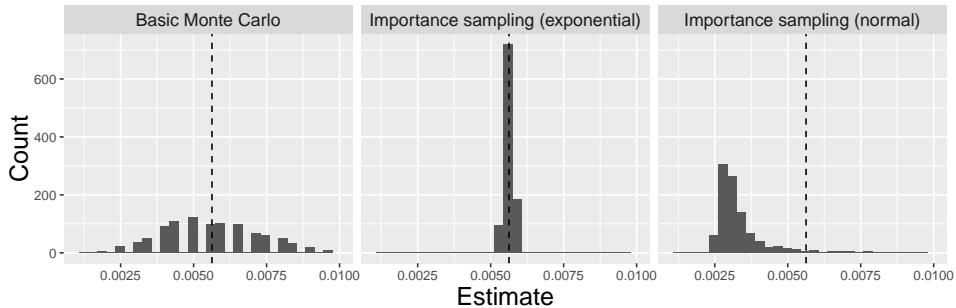


Figure 4: Distribution of estimates of the tail probability for an Ex-Gaussian distribution using basic Monte Carlo integration, importance sampling with a shifted exponential distribution, and importance sampling with a truncated Normal distribution. The dotted lines show the true value of the tail probability. Each estimator used $N = 2000$ particles and distributions were obtained by applying each estimator repeatedly for 1000 times.

235 heavier-tailed importance distribution which gives much better results.

236 *1.2. Efficiency*

237 As the previous example illustrates, some importance distributions are
 238 better than others. In the extreme case, a bad choice of importance dis-
 239 tribution can result in an estimator with an infinite variance. The optimal
 240 importance distribution q^* , in terms of minimizing the variance of the esti-
 241 mator, is

$$242 \quad q^*(y) = \frac{|f(y)|p(y)}{\int |f(y)|p(y) dy} \quad (3)$$

243 (Kahn and Marshall, 1953). For example, the optimal importance distri-
 244 bution in the previous example is a truncated ex-Gaussian. The optimal
 245 importance distribution is mainly of theoretical interest; to be able to use
 246 it, we would have to know the value of the integral $\int |f(z)|p(z) dz$, which is
 247 pretty much the quantity we want to estimate in the first place. Neverthe-
 248 less, we should aim to use an importance distribution which is as close to

249 this distribution as possible.

250 A more practical way to reduce the variance of the estimator is to nor-
251 malize the importance weights so they sum to 1 (Casella and Robert, 1998).
252 Using normalized weights

$$253 \quad W^{(i)} \triangleq \frac{w^{(i)}}{\sum_{j=1}^N w^{(j)}} \quad (4)$$

254 results in the “self-normalised” IS estimator

$$255 \quad E^{\text{ISn}} \triangleq \sum_{i=1}^N W^{(i)} f(y^{(i)}). \quad (5)$$

256 It should be noted that the estimator E^{ISn} is biased, but the bias is gen-
257 erally small, diminishes as the number of particles increases, and is often
258 offset by the gain in efficiency (the reduction in the variance of the estima-
259 tor).² Self-normalized weights are particularly useful in situations where the
260 distribution p is only known up to a normalizing constant. For instance, we
261 may be interested in the conditional distribution $p(y|x) = \frac{p(x,y)}{\int p(x,y) dy}$, but al-
262 though we can compute $p(x, y)$, the marginal distribution $p(x) = \int p(x, y) dy$
263 is intractable. In that case, we can still use importance sampling, as the
264 normalizing constant cancels out in the computation of the self-normalized

²One reason why the variance of the self-normalized IS estimate can be smaller is that the expected value of each weight is

$$\mathbb{E}_q[w(Y)] = \int \frac{p(y)}{q(y)} q(y) dy = \int p(y) dy = 1,$$

and the expected value of the sum of the weights is thus N . In practice, the summed importance weights will deviate from this value. Using self-normalizing weights ensures that the sum of the weights is always equal to 1 (not N , but we have accounted for this by removing the $1/N$ term in Equation 5), which removes one source of variance.

265 weights³. So, when using self-normalized weights, the target distribution
 266 only has to be known up to the normalizing constant.

267 **2. Sequential importance sampling and online Bayesian inference**

268 We often need to infer unknown parameters of a statistical model se-
 269 quentially after each new observation comes in. Such online inference is
 270 crucial in a wide range of situations, including adaptive design of exper-
 271 iments (e.g., Amzal, Bois, Parent, and Robert, 2006; Myung, Cavagnaro,
 272 and Pitt, 2013) and real-time fault diagnosis in nuclear power plants. From
 273 a Bayesian viewpoint, this means we want to compute a sequence of poste-
 274 rior distributions $p(\theta|y_1), p(\theta|y_{1:2}), \dots, p(\theta|y_{1:t})$, where $y_{1:t} = (y_1, y_2, \dots, y_t)$
 275 denotes a sequence of observations, and θ a vector of parameters. To ap-
 276 proximate such a posterior distribution $p(\theta|y_{1:t})$ with importance sampling,
 277 we need an importance distribution $q_t(\theta)$ to generate an importance sample
 278 of particles, and compute the importance weights

$$279 \quad w_t^{(i)} = \frac{p(\theta^{(i)}|y_{1:t})}{q_t(\theta^{(i)})}.$$

280 While we could generate a fresh importance sample at each point in time,
 281 this will usually increase the computational burden at each consecutive time
 282 point, as at each time we would have to browse the whole history of obser-
 283 vations to compute the importance weights. Moreover, when tracking an
 284 evolving latent state over time, we would also have to generate larger and

³This is shown as

$$\sum_{i=1}^N \frac{p(y^{(i)}|x) f(y^{(i)})}{q(y^{(i)})} = \sum_{i=1}^N \frac{\frac{1}{p(x)} \frac{p(x, y^{(i)})}{q(y^{(i)})} f(y^{(i)})}{\frac{1}{p(x)} \sum_{j=1}^N \frac{p(x, y^{(j)})}{q(y^{(j)})}} = \sum_{j=1}^N \frac{\frac{p(x, y^{(j)})}{q(y^{(j)})} f(y^{(j)})}{\sum_{j=1}^N \frac{p(x, y^{(j)})}{q(y^{(j)})}}$$

285 larger importance samples as, at each time point, we would have to sample
 286 the whole trajectory of the latent state thus far. For real-time applications,
 287 it is important to devise an algorithm with an approximately fixed computa-
 288 tional cost at each time point. Sequential importance sampling (SIS) serves
 289 this purpose. In addition, by using information from previous observations
 290 and samples, SIS can provide more efficient importance distributions than
 291 a straightforward application of IS.

292 A key idea in SIS is to compute the importance weights incrementally,
 293 by multiplying the importance weight at the previous time $t - 1$ by an incre-
 294 mental weight update $a_t^{(i)}$. It is always possible to formulate the importance
 295 weights in such a form by trivially rewriting the importance weights above
 296 as

$$\begin{aligned}
 297 \quad w_t^{(i)} &= \frac{p(\theta^{(i)}|y_{1:t})q_{t-1}(\theta^{(i)})}{p(\theta^{(i)}|y_{1:t-1})q_t(\theta^{(i)})} \frac{p(\theta^{(i)}|y_{1:t-1})}{q_{t-1}(\theta^{(i)})} \\
 298 \quad &= a_t^{(i)} w_{t-1}^{(i)},
 \end{aligned} \tag{6}$$

300 where we define the incremental weight update as

$$301 \quad a_t^{(i)} \triangleq \frac{p(\theta^{(i)}|y_{1:t})}{p(\theta^{(i)}|y_{1:t-1})} \times \frac{q_{t-1}(\theta^{(i)})}{q_t(\theta^{(i)})}. \tag{7}$$

302 This is of course not immediately helpful, as we still need to compute
 303 $p(\theta^{(i)}|y_{1:t})$ and $q_t(\theta^{(i)})$ in full at each time point. However, there are some im-
 304 portant cases where we can simplify the incremental weight update further.
 305 In this section, we will focus on one such case, involving on-line inference of
 306 time-invariant parameters. This will help to illustrate the basics of SIS and
 307 its main shortcoming. We will see that while sequential importance sampling
 308 is computationally efficient, allowing one to approximate the distributions
 309 of interest sequentially without having to revisit all previous observations
 310 or completely redraw the whole importance sample, the performance of SIS

311 degrades over time because after a large number of iterations, all but one
 312 particle will have negligible weight. After introducing resampling as a way
 313 to overcome this problem of “weight degeneracy”, we then return to a second
 314 important application of SIS, namely the inference of latent states in state-
 315 space models. Combining SIS with resampling provides us with a recipe for
 316 a particle filter which is highly effective for these applications.

317 *2.1. Sequential importance sampling for time-invariant parameters*

318 We will now focus on the problem of computing a sequence of poste-
 319 rior distributions $p(\theta|y_{1:t})$, $t = 1, \dots, T$ for a vector of unknown parameters
 320 θ . Assuming the observations are conditionally independent given the pa-
 321 rameters, we can express each of these posteriors through Bayes’ theorem
 322 as

323
$$p(\theta|y_{1:t}) = \frac{p(\theta) \prod_{i=1}^t p(y_i|\theta)}{p(y_{1:t})}.$$

324 In this case, the left-hand ratio in (7) simplifies to

325
$$\frac{p(\theta^{(i)}|y_{1:t})}{p(\theta^{(i)}|y_{1:t-1})} = p(y_t|\theta^{(i)})p(y_t|y_{1:t-1})$$

326 If we also use a single importance distribution $q_t(\theta) = q_{t-1}(\theta) = q(\theta)$ to ap-
 327 proximate each posterior, the right-hand ratio in (7) evaluates to $\frac{q_{t-1}(\theta^{(i)})}{q_t(\theta^{(i)})} =$
 328 1, so that the incremental weight update is simply $a_t^{(i)} = p(y_t|\theta^{(i)})p(y_t|y_{1:t-1})$.
 329 Using self-normalized importance weights, we can ignore the $p(y_t|y_{1:t-1})$
 330 term, resulting in a simple importance sampling scheme where we sequen-
 331 tially update the weights of an initial importance sample in light of each
 332 new observation:

333 **Algorithm 3.** SIS for time-invariant parameters

- 334 1. (Initialize) For $i = 1, \dots, N$, sample $\theta^{(i)} \sim q(\theta)$, and compute the
 335 normalized weights $W_0^{(i)} \propto \frac{p(\theta)}{q(\theta)}$ with $\sum_{j=1}^N W_0^{(j)} = 1$.

336 2. For $t = 1, \dots, t$:

337 (a) (Reweight) For $i = 1, \dots, N$, compute $W_t^{(i)} \propto p(y_t | \theta^{(i)}) W_{t-1}^{(i)}$,
338 with $\sum_{i=1}^N W_t^{(i)} = 1$.

339 (b) (Estimate) Compute the (self-normalized) SIS estimate

$$340 E_t^{\text{SISn}} = \sum_{i=1}^N W_t^{(i)} f(\theta^{(i)})$$

341 *2.1.1. Example: inferring the mean and variance of a Gaussian variable*

342 To illustrate how this algorithm works, we will apply it to sequentially
343 infer the (posterior) mean and variance of a random variable. The observa-
344 tions are assumed to be independent samples from a Gaussian distribution
345 with unknown mean μ and variance σ^2 , so the parameters are $\theta = (\mu, \sigma)$. As
346 prior distributions, we will use a Gaussian distribution for μ (with a mean
347 of 0 and a standard deviation of 10) and a uniform distribution for σ (in the
348 range between 0 and 50). As these priors are easy to draw from, we use them
349 also as importance distributions. We apply the algorithm to a total of 100
350 observations from a Gaussian distribution with mean $\mu = 5$ and standard
351 deviation $\sigma = 5$, using an importance sample of size $N = 200$ (note that
352 the sample size is kept low to enhance the clarity of the results; in real ap-
353 plications a larger sample size would be recommended). Figure 5 shows the
354 resulting estimates (posterior means) as well as the normalized importance
355 weights for each particle. We can see that the estimated posterior mean
356 of σ comes reasonably close to the true value as t increases. However, the
357 estimated posterior mean of μ converges to a value which is further off the
358 true value. The problem is that, as t increases, the weight of almost all the
359 particles becomes negligible. In the end, the posterior mean is effectively
360 estimated by a single particle $\theta^{(*)} = (\mu^{(*)}, \sigma^{(*)})$. While this single particle

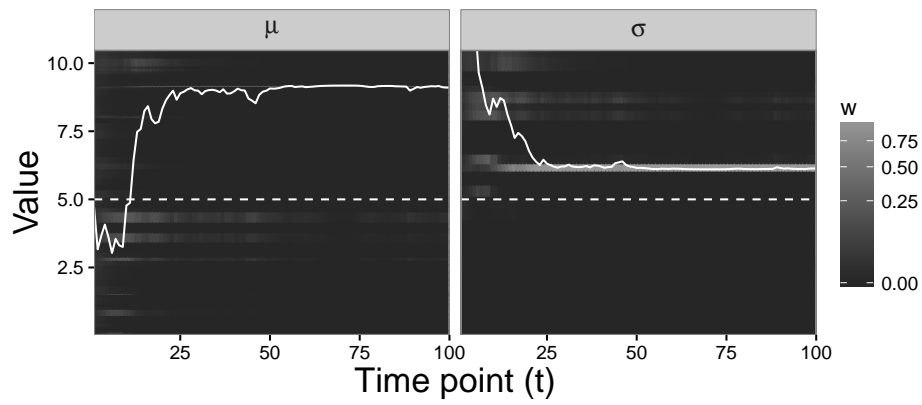


Figure 5: Online Bayesian inference of μ and σ of a Normal distribution. Solid white lines represent the posterior mean after each sequential observation and broken white lines the true values. Tiles in the background are centered around the particle values, such that each edge lies halfway between a particle and the immediately adjacent particle. The shade of each tile reflects the normalized weight of the corresponding particle. Initially, weights are almost uniformly distributed over the particles, but in the end (at time $t = 100$) only a single particle has non-negligible weight.

361 is, in some sense, the best one in the set, it does not have to be close to
362 the true values. In this run of the algorithm, $\sigma^{(*)}$ was quite close to the
363 true value, but $\mu^{(*)}$ was relatively far off from the true value. Indeed, there
364 were other particles with values $\mu^{(i)}$ which were closer to the true value.
365 However, for these particles, $\sigma^{(i)}$ was further off the true value, such that
366 taken together as a pair of parameter values, particle $\theta^{(*)}$ was better than
367 any of the other ones. The problem that the weights of almost all particles
368 approach 0 is referred to as weight degeneracy and a driving force behind it
369 is that the importance distribution becomes less and less efficient over time.
370 While the posterior distribution at first is close to the prior distribution
371 that was used as importance distribution, as more observations come in,
372 the posterior distribution becomes more and more peaked around the true
373 parameter values. The prior distribution is then too dispersed compared to
374 the posterior distribution and far from optimal as importance distribution.

375 As illustrated in Figure 2, we can think of the set of particles as defining
376 a random and irregularly-spaced grid and the SIS algorithm as approximat-
377 ing posterior distributions by computing the posterior probability of each
378 grid point. The algorithm becomes less and less efficient because after sam-
379 pling the particles, the grid is fixed. To make the algorithm more efficient,
380 we should adapt the grid points to each posterior distribution we want to
381 approximate. This is precisely what SMC algorithms do by resampling from
382 the particles (grid points), replicating particles with high and eliminating
383 those with low posterior probability at t . Insofar as the posterior probability
384 of particles at time $t+1$ is not wildly different from the posterior probability
385 at time t , this will thus provide useful grid points for time $t+1$. However,
386 resampling provides exact replicates of useful particles, and using the same
387 grid point multiple times does not increase the precision of the estimate.

388 After resampling, the particles are therefore “jittered” to rejuvenate the set,
 389 increasing the number of unique grid points and hence the precision of the
 390 approximation. Such jittering is natural for time-varying parameters. For
 391 instance, if instead of a fixed mean μ , we assume the mean μ_t changes from
 392 time to time according to a transition distribution $p(\mu_t|\mu_{t-1})$, we can use
 393 this transition distribution to generate samples of the current mean based on
 394 the samples of the previous mean. When the parameters are time-invariant,
 395 there is no immediately obvious way to “jitter” the particles after resam-
 396 pling. Some solutions have been proposed and we return to the problem
 397 of estimating static parameters with SMC in section 5.2. We will now first
 398 describe how resampling can be combined with (sequential) importance sam-
 399 pling, before turning to particle filters, which iterate sequential importance
 400 sampling and resampling steps to allow for flexible and efficient approxima-
 401 tions to posterior distributions of latent states in general state-space models.

402 3. Resampling

403 Due to the problem of weight degeneracy, after running an SIS algorithm
 404 for a large number of iterations (time points), all but one particle will have
 405 negligible weight. Clearly, this is not a good situation, as we effectively ap-
 406 proximate a distribution with a single particle. Moreover, this single particle
 407 does not even have to be in a region of high probability. A particle can have
 408 such a large relative weight that it will take very long for it to reduce, even
 409 though the target distribution has “moved on”. A useful measure to detect
 410 weight degeneracy is the *effective sample size* (Liu, 2001, p. 35-36), defined
 411 as

$$412 N_{\text{eff}} \triangleq \frac{1}{\sum_{i=1}^N (W^{(i)})^2} \quad (8)$$

413 This measure varies between 1 (all but one particle have weight 0) and N
414 (all particles have equal weight). Thus, the lower the effective sample size,
415 the stronger the weight degeneracy.

416 To counter the problem of weight degeneracy, SMC algorithms include
417 a resampling step, in which particles are sampled with replacement from
418 the set of all particles, with a probability that depends on the importance
419 weights. The main idea is to replicate particles with large weights and
420 eliminate those with small weights, as the latter have little effect on the
421 estimates anyway. The simplest sampling scheme is multinomial sampling,
422 which draws N samples from a multinomial distribution over the particle
423 indices $i = 1, \dots, N$, with probabilities $p(i) = W^{(i)}$. After resampling,
424 the weights are set to $W^{(i)} = 1/N$, because, roughly put, we have already
425 used the information in the weights to resample the set of particles. It is
426 straightforward to show that resampling does not change the expected value
427 of the estimator (5).

428 In addition to countering weight degeneracy, a further benefit of resam-
429 pling is that while the SIS samples themselves are not distributed according
430 to the target distribution p (they are distributed according to the instru-
431 mental distribution q and to approximate p we need to use the importance
432 weights), the resampled values are (approximately) distributed according
433 to p . A drawback of resampling is that it increases the variance of the es-
434 timator. To reduce this effect of resampling, alternatives to multinomial
435 resampling have been proposed with smaller variance.

436 The idea behind *residual resampling* (Liu and Chen, 1998) is to use a
437 deterministic approach as much as possible, and then use random resam-
438 pling for the remainder. To preserve the expected value of the estimator,
439 the expected number of replications of each particle i should be $NW^{(i)}$.

440 This is generally not an integer, and hence we can't use these expectations
 441 directly to generate the desired number of replications. Residual resampling
 442 takes the integer part of each $NW^{(i)}$ term and replicates each particle de-
 443 terministically according to that number. The remaining particles are then
 444 generated through multinomial resampling from a distribution determined
 445 by the non-integer parts of each $NW^{(i)}$ term.

446 *Stratified resampling* (Carpenter, Clifford, and Fearnhead, 1999) is an
 447 other scheme which results in partly deterministic replication of particles.
 448 As the name suggests, it is based on the principles of stratified sampling used
 449 in survey research. Practically, the method consists of using the weights to
 450 form an “empirical” cumulative distribution over the particles. This distri-
 451 bution is then split into N equally sized strata, and a single draw is taken
 452 from each stratum.

453 The most popular resampling scheme, *systematic resampling* (Kitagawa,
 454 1996), is based on the same intuition as stratified resampling, but reduces the
 455 Monte Carlo variance further by using a single random number, rather than
 456 a different random number, to draw from each stratum. Letting $\{\theta_t, W_t^{(i)}\}$
 457 represent the set of particles before resampling, and $\{\tilde{\theta}_t^{(i)}, \tilde{W}_t^{(i)}\}$ the set
 458 of particles after resampling, systematic resampling can be summarized as
 459 follows:

460 **Algorithm 4.** Systematic resampling

- 461 1. Draw $u \sim \text{Unif}(0, 1/N)$.
- 462 2. Define $U^i = (i - 1)/N + u$, $i = 1, \dots, N$.
- 463 3. For $i = 1, \dots, N$, find r such that $\sum_{k=1}^{r-1} W_t^{(k)} \leq U^i < \sum_{k=1}^r W_t^{(k)}$ and
 464 set $j(i) = r$.
- 465 4. For $i = 1, \dots, N$, set $\tilde{\theta}_t^{(i)} = \theta_t^{(j(i))}$ and $\tilde{W}_t^{(i)} = 1/N$.

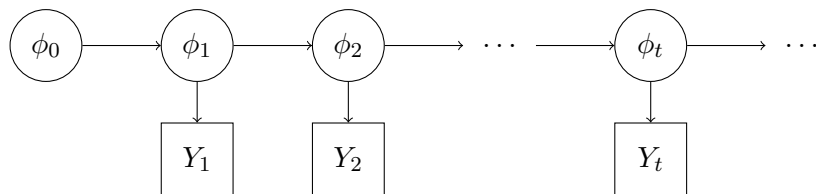


Figure 6: Schematic representation of a state-space model. Each observation Y_t at time t depends only on the current state ϕ_t . Each consecutive state ϕ_t depends only on the previous state ϕ_{t-1} .

466 Systematic resampling is simple to implement and generally performs very
 467 well in practice (Doucet and Johansen, 2011). However, in contrast to resid-
 468 ual and stratified resampling, it is not guaranteed to outperform multino-
 469 mial resampling (see Douc, Cappé, and Moulines, 2005, for this point and
 470 a thorough comparison of the theoretical properties of different resampling
 471 schemes).

472 4. Particle filters: SMC for state-space models

473 In filtering problems, we are interested in tracking the latent states ϕ_t of a
 474 stochastic process as each observation y_t comes in. In a Bayesian framework,
 475 we do this by computing a sequence of posterior distributions $p(\phi_1|y_1), \dots,$
 476 $p(\phi_t|y_{1:t})$. These posterior distributions of the current state ϕ_t , given all the
 477 observations thus far, are also called *filtering distributions*.

478 4.1. State-space models

479 State-space models are an important class of models to describe time-
 480 series of observations. State-space models describe an observable time series
 481 $y_{1:t}$ through a time-series of latent or hidden states $\phi_{0:t}$. In state-space
 482 models, we make two important assumptions about the relation between

483 states and observations. A graphical representation of a state-space model
 484 in the form of a Bayesian network is given in Figure 6. Firstly, we assume
 485 that each observation y_t depends solely on the current state ϕ_t , such that
 486 the observations are conditionally independent given the states ϕ_t :

$$487 \quad p(y_{1:T}|\phi_{0:T}) = \prod_{t=1}^T p(y_t|\phi_t),$$

488 Secondly, we assume that the hidden states change over time according to a
 489 first-order Markov process, such that the current state depends only on the
 490 state at the immediately preceding time point:

$$491 \quad p(\phi_{0:T}) = p(\phi_0) \prod_{t=1}^T p(\phi_t|\phi_{t-1}).$$

492 Given these two assumptions, we can write the posterior distribution over
 493 the hidden states as

$$494 \quad p(\phi_{0:T}|y_{1:T}) = \frac{p(\phi_0) \prod_{t=1}^T p(y_t|\phi_t)p(\phi_t|\phi_{t-1})}{p(y_{1:T})}.$$

495 Moreover, we can compute the posteriors recursively as

$$496 \quad p(\phi_{0:t}|y_{1:t}) = \frac{p(y_t|\phi_t)p(\phi_t|\phi_{t-1})}{p(y_t|y_{1:t-1})} p(\phi_{0:t-1}|y_{1:t-1}) \quad (9)$$

497 where

$$498 \quad p(y_t|y_{1:t-1}) = \iint p(y_t|\phi_t)p(\phi_t|\phi_{t-1})p(\phi_{t-1}|y_{1:t-1}) \mathrm{d}\phi_{t-1}\mathrm{d}\phi_t$$

499 Models with this structure are also known as hidden Markov models (e.g.,
 500 Visser, 2011).

501 4.2. SIS for state-space models

502 We can view the problem of estimating the hidden states $\phi_{0:t}$ as esti-
 503 mating a vector of parameters θ which increases in dimension at each time

504 point t , such that, at time t , we estimate $\theta = \phi_{0:t}$, and at time $t + 1$ we add
505 another dimension to the parameter vector to estimate $\theta = \phi_{0:t+1}$. Using
506 IS, we could draw a new importance sample $\{\phi_{0:t}^{(i)}\}$ at each time t , but this
507 would result in an increase in computational burden over time, which we
508 would like to avoid in real-time applications. Using sequential importance
509 sampling, we can incrementally build up the importance sample, starting at
510 time $t = 0$ with a sample $\{\phi_0^{(i)}\}$, then sampling values $\phi_1^{(i)}$ at time 1 con-
511 ditional on the sample at time 0, then adding sampled values $\phi_2^{(i)}$ at time
512 2 conditional on the sample at time 1, etc. Formally, this means we define
513 the importance distribution at time t as

$$514 \quad q_t(\phi_{0:t}) = q_t(\phi_t|\phi_{0:t-1})q_{t-1}(\phi_{0:t-1}). \quad (10)$$

515 Using this conditional importance distribution, and noting that $q_{t-1}(\phi_{0:t}) =$
516 $q_{t-1}(\phi_{0:t-1})$, the right-hand ratio in (7) simplifies to $\frac{q_{t-1}(\phi_{0:t})}{q_t(\phi_{0:t})} = \frac{1}{q_t(\phi_t|\phi_{0:t-1})}$.
517 Combining this with Equation 9, we can write the incremental weight update
518 as

$$519 \quad a_t^{(i)} = \frac{p(y_t|\phi_t^{(i)})p(\phi_t^{(i)}|\phi_{t-1}^{(i)})}{p(y_t|y_{1:t-1})q_t(\phi_t^{(i)}|\phi_{0:t-1}^{(i)})}.$$

520 Using normalized importance weights, we can ignore the $p(y_t|y_{1:t-1})$ term,
521 which is often difficult to compute. As when using SIS to sequentially es-
522 timate time-invariant parameters, we now have an algorithm of (approx-
523 imately) constant computational cost. At each time t , we add a new dimen-
524 sion to our particles by drawing values $\phi_t^{(i)}$ from a conditional importance
525 distribution, and we update the weights without the need to revisit all the
526 previous observations and hidden states.

527 As usual, the choice of importance distribution $q_t(\phi_t|\phi_{0:t-1})$ should be
528 carefully considered, as a poor choice can result in an estimator with high

529 variance. An equivalent formulation of the incremental weight update is

$$530 \quad a_t^{(i)} = \frac{p(\phi_t^{(i)} | y_t, \phi_{t-1}^{(i)}) p(y_t | \phi_{t-1}^{(i)})}{p(y_t | y_{1:t-1}) q_t(\phi_t^{(i)} | \phi_{0:t-1}^{(i)})},$$

531 which is generally more involved to compute, but indicates that the op-
 532 timal importance distribution (in terms of minimizing the variance of the
 533 importance weights) is

$$534 \quad q_t^*(\phi_t | \phi_{0:t-1}) = p(\phi_t | y_t, \phi_{t-1}).$$

535 The optimal importance distribution is again mostly of theoretical interest,
 536 but when possible, an importance distribution should be used which matches
 537 it as closely as possible.

538 Unfortunately, even when the optimal importance distribution can be
 539 used, SIS for state-space models will suffer from the same weight degeneracy
 540 problem we observed when estimating time-invariant parameters. Suppose
 541 that at time t , we have a “perfect” sample $\phi_{0:t}^{(i)} \sim p(\phi_{0:t} | y_{1:t})$, i.e., the parti-
 542 cles are distributed according to the target distribution and the normalized
 543 importance weights are all $W^{(i)} = 1/N$. Moving to the next time point, we
 544 can sample the new state $\phi_{t+1}^{(i)} \sim p(\phi_{t+1} | y_{t+1}, \phi_t^{(i)})$, but without redrawing
 545 $\phi_{0:t}^{(i)}$, the resulting particles $\phi_{0:t+1}^{(i)}$ will not be distributed according to the
 546 current target distribution $p(\phi_{0:t+1} | y_{1:t+1})$. Thus, in a sequential algorithm
 547 where we keep the values $\phi_{0:t}^{(i)}$, it will not be possible to generate a perfect
 548 sample at time $t + 1$. To do this, the samples $\phi_{0:t}^{(i)}$ would already have to be
 549 distributed according to $p(\phi_{0:t} | y_{1:t+1})$, and not according to $p(\phi_{0:t} | y_{1:t})$, the
 550 target distribution at time t . These two distributions are generally differ-
 551 ent, because the new observation y_{t+1} provides information about the likely
 552 values of $\phi_{0:t}$ that was not available at the time of drawing $\phi_{0:t}^{(i)}$. Hence,

553 repeated application of SIS necessarily introduces variability in the impor-
 554 tance weights, which builds up over time, increasing the variance of the
 555 estimator and ultimately resulting in weight degeneracy.

556 *4.3. A generic particle filter*

557 To counter weight degeneracy, SMC methods combine SIS with resam-
 558 pling, replicating particles with high weights and eliminating those with low
 559 weights. After resampling, the particles have uniform weights and are (ap-
 560 proximately) distributed according to the target distribution $p(\phi_{0:t}|y_{1:t})$. At
 561 the next time point, a new dimension ϕ_{t+1} is added to the particles by draw-
 562 ing from a conditional importance distribution $q_{t+1}(\phi_{t+1}^{(i)}|\phi_{0:t}^{(i)})$. Resampling
 563 will have focussed the set $\{\phi_{0:t}^{(i)}\}$ on useful “grid points” and while this set
 564 contains exact copies of particles, the values of the new dimension $\phi_{t+1}^{(i)}$ will
 565 be jittered and hence provide an adapted and useful set of grid points to es-
 566 timate ϕ_{t+1} . As estimation of a current state ϕ_t is generally of more interest
 567 than estimating the whole path $\phi_{0:t}$, the fact that we now have a multidim-
 568 ensional grid where only the last dimension is adapted (as the dimensions
 569 reflecting earlier states contain exact replicates) is of little concern. If we
 570 are interested in estimating the whole path $\phi_{0:t}$, we would need to jitter the
 571 grid points on all dimensions. We return to this issue in Section 5.1.

572 A generic particle filter (see Figure 1) to approximate a sequence of
 573 posterior distributions $p(\phi_{0:1}|y_1)$, $p(\phi_{0:2}|y_{1:2})$, \dots , $p(\phi_{0:t}|y_{1:t})$, proceeds as
 574 follows:

575 **Algorithm 5.** A generic particle filter

- 576 1. (Initialize) For $i = 1, \dots, N$, sample $\tilde{\phi}_0^{(i)} \sim q(\phi_0)$ and compute the
 577 normalized importance weights $\tilde{W}_0^{(i)} \propto \frac{p(\tilde{\phi}_0^{(i)})}{q(\tilde{\phi}_0^{(i)})}$ with $\sum_{i=1}^N \tilde{W}_0^{(i)} = 1$.

578 2. For $t = 1, \dots, T$:

579 (a) (Propagate) For $i = 1, \dots, N$, sample $\phi_t^{(i)} \sim q_t(\phi_t | \tilde{\phi}_{0:t-1}^{(i)})$, and
580 add this new dimension to the particles, setting $\phi_{0:t}^{(i)} = (\tilde{\phi}_{0:t-1}^{(i)}, \phi_t^{(i)})$.

581 (b) (Reweight) For $i = 1, \dots, N$, compute normalized weights

$$582 \quad W_t^{(i)} \propto \frac{p(y_t | \phi_t^{(i)}) p(\phi_t^{(i)} | \phi_{t-1}^{(i)})}{q_t(\phi_t^{(i)} | \phi_{0:t-1}^{(i)})} \tilde{W}_{t-1}^{(i)}$$

583 with $\sum_{i=1}^N W_t^{(i)} = 1$.

584 (c) (Estimate) Compute the required estimate

$$585 \quad E_t^{\text{PFn}} = \sum_{i=1}^N f(\phi_{0:t}^{(i)}) W_t^{(i)}.$$

586 (d) (Resample) If $N_{\text{eff}} \leq cN$, resample $\{\tilde{\phi}_{0:t}^{(i)}\}$ with replacement from
587 $\{\phi_{0:t}^{(i)}\}$ using the normalized weights $W_t^{(i)}$ and set $\tilde{W}_t^{(i)} = 1/N$ to
588 obtain a set of equally weighted particles $\{\tilde{\phi}_t^{(i)}, \tilde{W}^{(i)} = 1/N\}$; else
589 set $\{\tilde{\phi}_{0:t}^{(i)}, \tilde{W}_t^{(i)}\} = \{\phi_{0:t}^{(i)}, W_t^{(i)}\}$.

590 In this generic particle filter, we allow for optional resampling, whenever
591 the effective sample size is smaller than or equal to a proportion c of the
592 number of particles used N . While particle filters often set $c = 1$, so that
593 resampling is done on every step, choosing a different value can be benefi-
594 cial, as resampling introduces additional variance in the estimates. If the
595 importance weights show little degeneracy, then this additional variance is
596 unnecessary. Therefore, setting $c = .5$, which is another common value, we
597 only resample when there is sufficient evidence for weight degeneracy.

598 At each iteration, we effectively have two particle approximations, the
599 set $\{\phi_{0:t}^{(i)}, W_t^{(i)}\}$ before resampling, and the set $\{\tilde{\phi}_{0:t}^{(i)}, \tilde{W}^{(i)}\}$ after resampling.

600 While both provide unbiased estimates, the estimator before resampling

601 generally has lower variance. It should also be noted that this particle filter
602 provides a weighted sample of state *sequences* $\phi_{0:t}^{(i)} = (\phi_0^{(i)}, \phi_1^{(i)}, \dots, \phi_t^{(i)})$,
603 approximating a posterior distribution over state sequences $p(\phi_{0:t}|y_{1:t})$. The
604 estimator E_t^{PFn} is also defined over these sequences and not a single state
605 ϕ_t . In filtering problems, we are generally only interested in estimating the
606 current state, not the whole path $\phi_{0:t}$. As the posterior distribution over a
607 single state is a marginal distribution of the joint posterior distribution over
608 all states, we can write the required expected value as

$$\begin{aligned}
609 \quad \mathbb{E}_p[f(\phi_t)] &= \int f(\phi_t)p(\phi_t|y_{1:t}) d\phi_t \\
610 &= \iint f(\phi_t)p(\phi_t, \phi_{0:t-1}|y_{1:t}) d\phi_t d\phi_{0:t-1}. \\
611
\end{aligned}$$

612 This means that we can effectively ignore the previous states, and use the
613 estimator

$$614 \quad E^{\text{PFn}} = \sum_{i=1}^n W_t^{(i)} f(\phi_t^{(i)}).$$

615 Several variants of this generic particle filter can be found in the litera-
616 ture. In the bootstrap filter (Gordon, Salmond, and Smith, 1993), the state
617 transition distribution is used as the conditional importance distribution, i.e.
618 $q_t(\phi_t|\phi_{0:t-1}^{(i)}) = p(\phi_t|\phi_{t-1}^{(i)})$. This usually makes the propagate step simple to
619 implement and also simplifies the reweighting step, as the weights can now
620 be computed as $W_t^{(i)} \propto p(y_t|\phi_t^{(i)})\tilde{W}_{t-1}^{(i)}$. The auxiliary particle filter, intro-
621 duced in Pitt and Shephard (1999) and later improved by Carpenter et al.
622 (1999), effectively switches the resampling and propagate steps, resulting in
623 a larger number of distinct particles to approximate the target. The aux-
624 iliary particle filter can be implemented as a variant of the generic particle
625 filter by adapting the importance weights to incorporate information from
626 the observation at the next time point; the subsequent resampling step is

627 then based on information from the next time point, and thus particles will
628 be resampled which are likely to be useful for predicting this time point. For
629 more information, see, e.g., Doucet and Johansen (2011) and Whiteley and
630 Johansen (2011).

631 Chopin (2004) shows how a general class of SMC algorithms, including
632 the generic particle filter, satisfy a Central Limit Theorem, such that, as the
633 number of particles approaches infinity, SMC estimates follow a Gaussian
634 distribution centered around the true value. For other convergence results
635 and proofs, see, e.g., Del Moral (2013), Douc and Moulines (2008), and
636 Whiteley (2013).

637 4.4. Example: A particle filter for a simple Gaussian process

638 We will illustrate SIS with an example of a latent Gaussian process with
639 noisy observations. Suppose there is a latent variable ϕ which moves in
640 discrete time according to a random walk

$$641 \quad \phi_{t+1} = \phi_t + \xi_t \quad \xi_t \sim N(0, \sigma_\xi^2), \quad (11)$$

642 where the initial distribution at $t = 0$ is given as

$$643 \quad \phi_0 \sim N(\mu_0, \sigma_0^2). \quad (12)$$

644 The value of the latent variable can only be inferred from noisy observations
645 Y_t that depend on the latent process through

$$646 \quad Y_t = \phi_t + \epsilon_t \quad \epsilon_t \sim N(0, \sigma_\epsilon^2). \quad (13)$$

647 This model is a relatively simple state-space model with

$$648 \quad p(y_t|\phi_t) = N(\phi_t, \sigma_\epsilon^2)$$

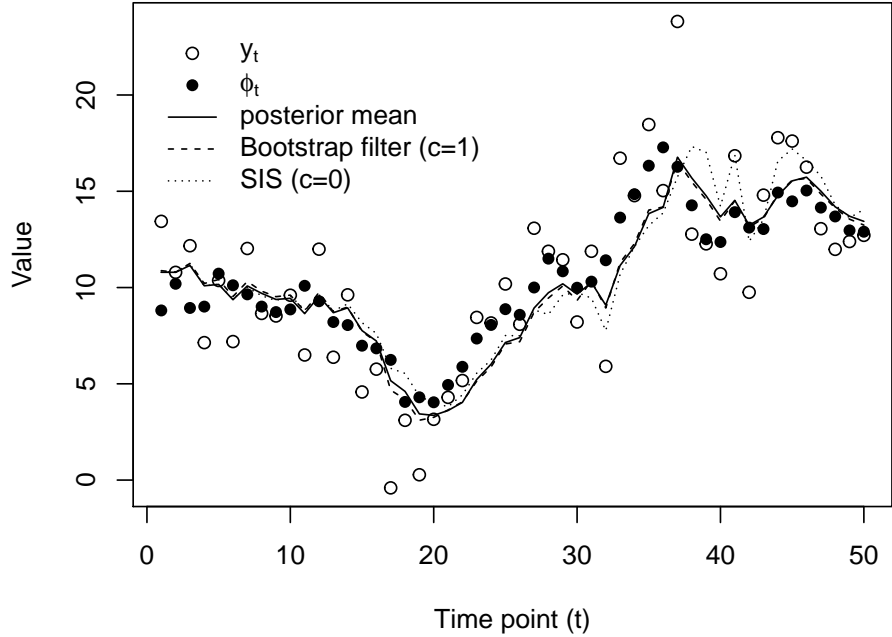


Figure 7: Example of the one-dimensional Gaussian latent process with $\mu_0 = 10$, $\sigma_0^2 = 2$, $\sigma_\xi^2 = 1$, and $\sigma_\epsilon^2 = 10$. Dots show the observations y_t and latent states ϕ_t . Lines show the true posterior means and particle filter estimates of the posterior means.

649 and

$$650 \quad p(\phi_t | \phi_{t-1}) = N(\phi_{t-1}, \sigma_\xi^2).$$

651 Figure 7 contains example data from this process.

652 Suppose that the observations y_t are made sequentially, and after each
 653 new observation, we wish to infer the value of the underlying latent variable
 654 ϕ_t . The distributions of interest are thus $p(\phi_1 | y_1)$, $p(\phi_2 | y_{1:2})$, \dots , $p(\phi_t | y_{1:t})$.
 655 As the process is linear and Gaussian, the distributions can be computed

656 analytically by the Kalman filter (Kalman, 1960; Kalman and Bucy, 1961).
 657 However, approximating the posterior means by a particle filter will illus-
 658 trate some key features of the algorithm and the availability of analytical
 659 estimates offers a useful standard to evaluate its quality.

660 We will use a bootstrap filter, using as conditional importance distribu-
 661 tions the transition distribution of the latent process, i.e. $q_t(\phi_t|\phi_{0:t-1}) =$
 662 $p(\phi_t|\phi_{t-1})$. The self-normalized weights are then easily computed as $W_t \propto$
 663 $p(y_t|\phi_t^{(i)})\tilde{W}_{t-1}^{(i)}$. We also set $c = 1$, so that we resample at each iteration,
 664 using systematic resampling. Figure 7 contains the resulting estimates E_t^{PFn}
 665 of the posterior means as well as the analytical (true) values computed with
 666 a Kalman filter. As can be seen there, the estimates are very close to the
 667 analytical posterior means. For comparison, if we run the filter with $c = 0$
 668 (so that we never resample, turning it into a straightforward SIS algorithm),
 669 we see that while the estimates are quite good initially, at later time points,
 670 the deviation between the estimated and actual posterior means increases.
 671 Again, this is due to weight degeneracy, which is countered by the resam-
 672 pling step in the particle filter. The effect of resampling can be clearly seen
 673 in Figure 8, which depicts the variation in the estimates when the algorithms
 674 are applied repeatedly to the same data. While there clearly is an increase
 675 in the variance of SIS over time, the estimates of the bootstrap filter remain
 676 close to the analytical values. For comparison, we also plot the results of an
 677 SIS and particle filter algorithm with the optimal importance distribution
 678 $q_t(\phi_t|\phi_{0:t-1}) = p(\phi_t|y_t, \phi_{t-1})$. The increase in efficiency due to the optimal
 679 importance distribution is clearly seen in the case of SIS. While still present,
 680 the difference between the two particle filters appears less marked.

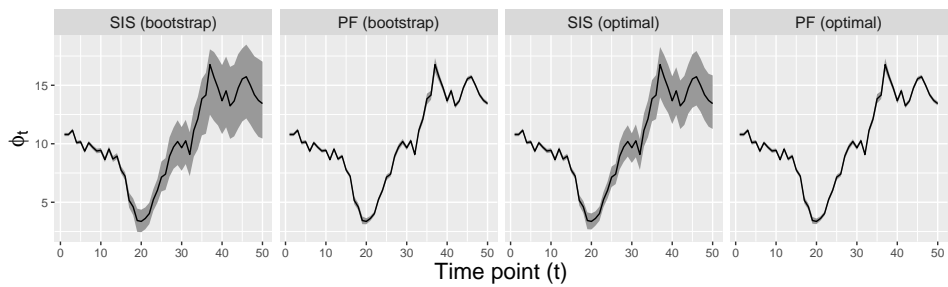


Figure 8: A comparison of SIS (no resampling) and particle filters (with resampling) shows that resampling clearly improves the accuracy of estimation. Each panel shows the true posterior mean (solid line) and 95% interpercentile region of estimated posterior means (shaded region) when applying the algorithms to the data of Figure 7. SIS (bootstrap) is the bootstrap filter with $c = 0$ (no resampling), PF (bootstrap) is the bootstrap filter with $c = 1$ (always resampling), and SIS (optimal) and PF (optimal) are similar to these but use the optimal importance distribution. Each algorithm uses $N = 500$ particles and the interpercentile regions were computed by running each algorithm on the same data for 2000 replications.

681 **5. Further issues and extensions**

682 While particle filters generally work well to approximate the filtering
683 distributions of latent states in general (non-linear and/or non-Gaussian)
684 state-space models, the application of SMC beyond this domain requires
685 further consideration. Here we briefly highlight some issues and solutions
686 for using SMC to approximate the posterior distributions over whole state
687 trajectories and time-invariant or static parameters. We end with a brief
688 discussion of how to combine sampling with analytical integration in order
689 to increase the efficiency of SMC.

690 *5.1. Sample impoverishment and particle smoothing*

691 Although resampling counters the problem of weight degeneracy, it in-
692 troduces a new problem of “sample impoverishment”. By replicating and
693 removing particles, resampling reduces the total number of unique values
694 present in the set of particles. This is no major issue in filtering, where we
695 are interested in estimating the current state ϕ_t . The particle values $\phi_t^{(i)}$
696 used for this are “jittered” in the propagate step, and estimation proceeds
697 before resampling affects the number of unique values of this state in the set
698 of particles. However, resampling does reduce the unique values of $\phi_{0:t-1}^{(i)}$
699 reflecting the states at earlier time points, and over time this problem be-
700 comes more and more severe for the initial states. When the algorithm has
701 run for a large number of time points, it can be the case that all the parti-
702 cles have the same value for ϕ_1 . This sample impoverishment can severely
703 affect the approximation of the so called *smoothing* distributions $p(\phi_t|y_{1:T})$,
704 $1 \leq t \leq T$.

705 To provide better approximations of smoothing distributions, several
706 alternatives to basic SMC have been proposed. A simple procedure is to use

707 a *fixed lag* approximation (Kitagawa and Sato, 2001). This approach relies
 708 on the exponential forgetting properties of many state-space models, such
 709 that

$$710 \quad p(\phi_{0:t}|y_{1:T}) \approx p(\phi_{0:t}|y_{1:(t+\Delta)}) \quad (14)$$

711 for a certain integer $0 < \Delta < T - t$; that is, observations after time $t + \Delta$
 712 provide no additional information about $\phi_{0:t}$. If this is the case, we do not
 713 need to update the estimates of $\phi_{0:t}$ after time $t + \Delta$. For the SMC algorithm,
 714 that means we do not need to update (e.g., resample) the particle values $\phi_{0:t}^{(i)}$
 715 then. Generally, we do not know Δ , and hence have to choose a value D
 716 which may be smaller or larger than Δ . If $D > \Delta$, we have not reduced the
 717 degeneracy problem as much as we could. If $D < \Delta$, then $p(\phi_{0:t}|y_{1:(t+D)})$ is
 718 a poor approximation of $p(\phi_{0:t}|y_{1:T})$.

719 A better, but computationally more expensive option is to store the par-
 720 ticle approximations of the filtering distributions (i.e., the particle values
 721 $\phi_t^{(i)}$ and weights $W_t^{(i)}$ approximating $p(\phi_t|y_{1:t})$) and then reweight these us-
 722 ing information from observations $y_{(t+1):T}$ to obtain an approximation of
 723 $p(\phi_t|y_{1:T})$. Particle variants of the Forward Filtering Backwards Smoothing
 724 and Forward Filtering Backwards Sampling algorithms have been proposed
 725 for this purpose (see e.g., Douc, Garivier, Moulines, and Olsson, 2011). One
 726 issue is that as these methods reweight or resample the particle values used
 727 to approximate the filtering distributions, but do not generate new particle
 728 values, they can be expected to perform poorly when the smoothing distri-
 729 butions differ substantially from the filtering distributions. An alternative
 730 approach, which can be expected to perform better in these situations, is
 731 the two-filter formulation of Briers, Doucet, and Maskell (2010).

732 5.2. *Inferring time-invariant (static) parameters*

733 While particle filtering works generally well to estimate latent states
 734 ϕ_t , which can be viewed as time-varying parameters, estimation of time-
 735 invariant or static parameters θ is more problematic. For instance, consider
 736 the simple Gaussian process defined in (11) - (13), and suppose both σ_ξ and
 737 σ_ϵ are unknown. The inference problem is then to approximate

$$738 \quad p(\phi_{0:t}, \theta | y_{1:t}) \propto p(\phi_{0:t} | \theta, y_{1:t}) p(\theta | y_{1:t})$$

739 where $\theta = (\sigma_\xi, \sigma_\epsilon)$. The main problem for a particle filter approximation is
 740 again sample impoverishment: resampling will reduce the number of unique
 741 particle values that represent θ . For time-invariant parameters, there is no
 742 natural way to “jitter” the particles after resampling. One solution is to use
 743 an artificial dynamic process for the static parameters, e.g., drawing new
 744 particles from a (multivariate) Normal distribution centered around the old
 745 particle values

$$746 \quad \theta_{t+1}^{(i)} \sim N(\theta_t^{(i)}, \Sigma_\theta),$$

747 where Σ_θ is a covariance matrix, and subscript t indicates that the parameter
 748 values reflect the time t posterior and not that the parameters are actually
 749 time-varying. While this reduces the problem of sample impoverishment, the
 750 artificial process will inflate the variance of the posterior distributions. Liu
 751 and West (2001) propose to view the artificial dynamics as a form of kernel
 752 smoothing: drawing new particle values from a Normal distribution centered
 753 around the old particle values is akin to approximating the distribution of
 754 θ by a (multivariate) Gaussian kernel density:

$$755 \quad p(\theta | y_{1:t}) \approx \sum_{i=1}^N W^{(i)} N(\theta | \theta_t^{(i)}, \Sigma_\theta).$$

756 To reduce variance inflation, Liu and West (2001) suggest to use a form of
 757 shrinkage, shifting the kernel locations $\theta_t^{(i)}$ closer to their overall mean by
 758 a factor which ensures that the variance of the particles equals the actual
 759 posterior variance.

760 An alternative is to incorporate Markov chain Monte Carlo (MCMC)
 761 moves in the particle filter (Andrieu, Doucet, and Holenstein, 2010; Chopin,
 762 2002; Chopin, Jacob, and Papaspiliopoulos, 2013; Gilks and Berzuini, 2001).
 763 The idea is to rejuvenate the particle set by applying a Markov transition
 764 kernel with the correct invariant distribution as the target. As they leave the
 765 target distribution intact, inclusion of MCMC moves in a particle filter al-
 766 gorithm is generally allowed. For a recent comparison of various approaches
 767 to parameter estimation with SMC techniques, see Kantas, Doucet, Singh,
 768 Maciejowski, and Chopin (2015).

769 5.3. Rao-Blackwellized particle filters

770 There are models in which, conditional upon some parameters, the dis-
 771 tributions of the remaining parameters can be solved analytically. For in-
 772 stance, in the example of the Gaussian process, we could assume that the
 773 process can switch between periods of high and low volatility. This can be
 774 represented by assuming a second latent process $\omega_{1:T}$, where ω_t is a discrete
 775 latent state indicating low or high volatility, and letting the innovation vari-
 776 ance $\sigma_\xi(\omega_t)$ be a function of this discrete latent state. This is an example
 777 of a switching linear state-space model, which is analytically intractable.
 778 Writing the joint posterior as

$$779 \quad p(\phi_{0:t}, \omega_{1:t} | y_{1:t}) = p(\phi_{0:t} | \omega_{1:t}, y_{1:t}) p(\omega_{1:t} | y_{1:t})$$

780 and realizing that, conditional upon $\omega_{1:t}$, $\phi_{0:t}$ is a linear Gaussian state-
 781 space model, the Kalman filter can be used to analytically compute the

782 conditional distributions $p(\phi_{0:t}|\omega_{1:t}, y_{1:t})$. Hence, we only need to approxi-
783 mate $p(\omega_{1:t}|y_{1:t})$ through sampling (cf. Chen and Liu, 2000). Solving part
784 of the problem analytically reduces the variance in the importance weights,
785 and hence increases the reliability of the estimates (Chopin, 2004). The
786 main message here is that sampling should be avoided whenever possible.
787 The combination of sampling and analytical inference is also called Rao-
788 Blackwellisation (Casella and Robert, 1996).

789 **6. A particle filter to track changes in learning rate during prob-** 790 **abilistic category learning**

791 As a final example of SMC estimation, we use a particle filter to esti-
792 mate changes in learning rate in probabilistic category learning. In proba-
793 bilistic category learning tasks, people learn to assign objects to mutually
794 exclusive categories according to their features, which are noisy indicators
795 of category membership. Tracking the learning rate throughout a task is
796 theoretically interesting, as it reflects how people adapt their learning to
797 the volatility in the task. If the relation between features and category
798 membership is time-invariant, people should ideally show “error discount-
799 ing” (Craig, Lewandowsky, and Little, 2011; Speekenbrink, Channon, and
800 Shanks, 2008), where they accept an unavoidable level of error and stabilize
801 their classification strategy by slowly stopping to learn. On the other hand,
802 if the relation between features and category membership changes over time,
803 then people should continue to adapt their categorization strategy to these
804 changes (Behrens, Woolrich, Walton, and Rushworth, 2007; Speekenbrink
805 and Shanks, 2010). How quickly they adapt (i.e. their learning rate) should
806 ideally depend on the rate at which the feature-category relation changes

807 (i.e. the volatility in the task). When the volatility is unknown, this it-
808 self has to be inferred from experience, such that people effectively have to
809 “learn how (much) to learn”.

810 Previous investigations of dynamic changes in learning rate have either
811 estimated learning rates separately in consecutive blocks of trials (Behrens
812 et al., 2007), or assumed the changes in learning rate followed a predeter-
813 mined schedule (Craig et al., 2011; Speekenbrink et al., 2008). Using a parti-
814 cle filter, we can estimate the learning rate on a trial-by-trial basis without
815 making too restrictive assumptions. In this example, we use unpublished
816 data collected in 2005 by David A. Lagnado at University College London.
817 Nineteen participants (12 female, average age 25.84) performed the Weather
818 Prediction Task (WPT, Knowlton, Squire, and Gluck, 1994). In the WPT,
819 the objective is to predict the state of the weather (“fine” or “rainy”) on the
820 basis of four cues (tarot cards with different geometric patterns, which are
821 either present or absent). Two cues are predictive of fine weather and two
822 cues of rainy weather. The version of the WPT used here included a sudden
823 change, whereby from trial 101 until the final (200th) trial, cues that were
824 first predictive of fine weather become predictive of rainy weather, and vice
825 versa.

826 As in Speekenbrink et al. (2008), we will assume people learn an asso-
827 ciative weight, v_j , for each cue. Over trials, these weights are updated by
828 the delta-rule

$$829 \quad v_{j,t+1} = v_{j,t} + \eta_t (y_t - p_t) x_{j,t},$$

830 where $x_{j,t}$ is a binary variable reflecting whether cue j was present on trial
831 t , y_t is a binary variable reflecting the state of the weather, and p_t is the

832 predicted probability of the state of the weather:

$$833 \quad p_t = \frac{1}{1 + \exp\left(-\sum_{j=1}^4 v_{j,t}x_{j,t}\right)}$$

834 People’s categorization responses are also assumed to follow these predicted
835 probabilities.

836 Our interest is in the learning rate $\eta_t > 0$, which we will allow to vary
837 from trial to trial according to the following transition distribution:

$$838 \quad p(\eta_{t+1}|\eta_t) = TN(\eta_t, \sigma_\eta^2) \quad (15)$$

839 where TN is a normal distribution truncated below at 0 (as the learning
840 rate is positive). We also use a truncated normal distribution for the initial
841 learning rates:

$$842 \quad \eta_0 \sim TN(\mu_0, \sigma_0^2)$$

843 Estimating the learning rates η_t is a difficult problem. Each response is
844 a random variable drawn from a Bernoulli distribution with parameter p_t ,
845 which depends on the cues on trial t and the associative weights $v_{j,t}$. These
846 associative weights in turn depend on the starting weights $v_{j,1}$ (which we
847 fix to 0), the previous cues and states of the weather, and the sequence of
848 unknown learning rates $\eta_{1:t-1}$. While the delta rule is deterministic, un-
849 certainty about the learning rates induces uncertainty about the associative
850 weights. Unfortunately, as we only have a single binary response on each
851 trial, we obtain relatively little information to reduce the uncertainty about
852 the associative weights and with that about the learning rates which gave
853 rise to these weights. All in all, we can thus expect the estimated learning
854 rates to be somewhat noisy.

855 For each participant, we estimated the time-varying learning rates with
856 a bootstrap filter with $N = 2000$ particles and selective resampling when

857 $N_{\text{eff}} < .5N$, i.e. when the effective sample size was half the number of
858 particles. The hyper-parameters were $\mu_0 = 0.05$, $\sigma_0 = 0.795$ and $\sigma_\eta = 0.101$,
859 which were determined maximizing the likelihood of the responses for all
860 participants.

861 We expected learning rates to start relatively high at the beginning of
862 the task and then to gradually decrease due to error discounting. In re-
863 sponse to the abrupt change in task structure at trial 101, learning rates
864 were expected to increase again, possibly decreasing again thereafter. Fig-
865 ure 9 shows the estimated means and 5% and 95% quantiles of the posterior
866 (filtering) distributions $p(\eta_t | \mathbf{x}_{1:t+1}, y_{1:t}, r_{1:t+1})$. The results show that many
867 participants show an increase in learning rate after trial 101, reflecting the
868 change in task structure at that point. Thus, many participants appeared to
869 indeed adapt their learning to the volatility in the environment. While some
870 participants, such as S5, show the expected pattern with initially relatively
871 high learning rate which decreases until the change at trial 101, then in-
872 creasing and decreasing again thereafter, other participants, such as S14 do
873 not show slowed learning towards the end of the task. This might be due to
874 expecting more abrupt changes. There is quite some individual variability in
875 the dynamics of learning rate over time. The best performing participants
876 have relatively high learning rates and marked changes throughout the task.
877 The less well performing participants had relatively low learning rates, in-
878 dicative of a slow adaptation of their strategy to the task structure. While
879 the posterior distributions are generally wide, because the responses provide
880 limited information about the learning rates, the results were consistent over
881 multiple runs of the algorithm and deviations in the hyper-parameters.

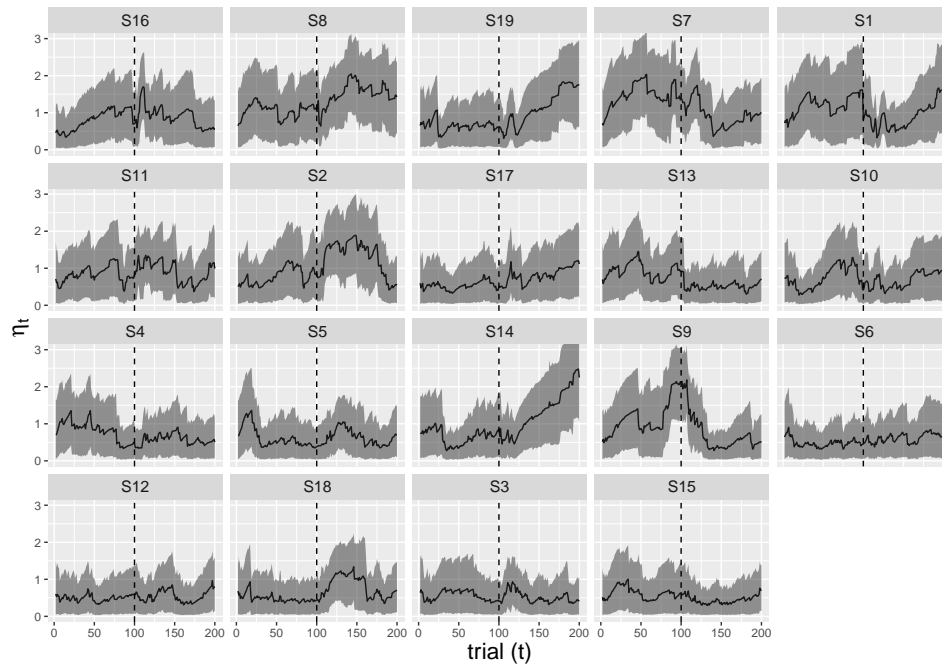


Figure 9: Particle filter estimates of time-varying learning rates. Solid lines represent the estimated mean (solid line) and shaded areas the 5%-95% interpercentile range of the posterior (filtering) distributions of η_t for each participant. Participants are ordered row-wise according to the expected performance of their predictions (the probability that their prediction was correct, rather than whether it was actually correct on a particular trial), with the best performing participant appearing in the top-left panel.

882 7. Conclusion

883 This tutorial introduced sequential Monte Carlo (SMC) estimation and,
884 in particular, particle filters. These techniques provide sampling-based ap-
885 proximations of a sequence of posterior distributions over parameter vec-
886 tors which increase in dimension, and allow inference in complex dynamic
887 statistical models. They rely on a combination of sequential importance
888 sampling and resampling steps. Sequential importance sampling provides
889 sets of weighted random samples (particles), while resampling reduces the
890 problem of weight degeneracy that plagues sequential importance sampling.
891 SMC has proven especially useful in filtering problems for general state-
892 space models, where the objective is to estimate the current value of a
893 latent state given all previous observations, but its use extends to other
894 problems including maximum likelihood estimation (e.g., Johansen, Doucet,
895 and Davy, 2008) and optimizing experimental designs (Amzal et al., 2006).
896 SMC is an active research field in statistics and machine learning and re-
897 cent developments have focussed on combining SMC with Markov chain
898 Monte Carlo (MCMC) techniques in order to provide efficient inference for
899 statistical models with both dynamic states and static parameters (e.g., An-
900 drieu et al., 2010; Chopin et al., 2013). Recent software to implement SMC
901 techniques include Biips (<http://alea.bordeaux.inria.fr/biips/>) and
902 LibBi (<http://libbi.org/>). All analyses in this tutorial were programmed
903 in the R language and all code and data are available in the supplementary
904 material and on the Open Science Framework (<http://osf.io/b6gsk/>).

905 **8. Acknowledgements**

906 This work was supported by the U.K. Economic and Social Research
907 Council (ESRC), grant RES-062-23-1511.

908 **9. References**

909 Amzal, B., Bois, F. Y., Parent, E., Robert, C. P., 2006. Bayesian-optimal
910 design via interacting particle systems. *Journal of the American Statistical*
911 *Association* 101, 773–785.

912 Andrieu, C., Doucet, A., Holenstein, R., 2010. Particle Markov chain Monte
913 Carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical*
914 *Methodology)* 72 (3), 269–342.

915 Behrens, T. E. J., Woolrich, M. W., Walton, M. E., Rushworth, M. F. S.,
916 Sep. 2007. Learning the value of information in an uncertain world. *Nature*
917 *neuroscience* 10 (9), 1214–21.

918 Briers, M., Doucet, A., Maskell, S., 2010. Smoothing algorithms for
919 state?space models. *Annals of the Institute of Statistical Mathematics*
920 62, 61–89.

921 Brown, S. D., Steyvers, M., 2009. Detecting and predicting changes. *Cognitive*
922 *Psychology* 58, 49–67.

923 Carpenter, J., Clifford, P., Fearnhead, P., 1999. Improved particle filter
924 for nonlinear problems. *IEE Proceedings-Radar, Sonar and Navigation*
925 146 (1), 2–7.

926 Casella, G., Robert, C. P., 1996. Rao-blackwellisation of sampling schemes.
927 *Biometrika* 83 (1), 81–94.

- 928 Casella, G., Robert, C. P., 1998. Post-processing accept-reject samples: Re-
929 cycling and rescaling. *Journal of Computational and Graphical Statistics*
930 7 (2), pp. 139–157.
- 931 Chen, R., Liu, J. S., 2000. Mixture kalman filters. *Journal of the Royal*
932 *Statistical Society: Series B (Statistical Methodology)* 62 (3), 493–508.
- 933 Chopin, N., 2002. A sequential particle filter method for static models.
934 *Biometrika* 89 (3), 539–552.
- 935 Chopin, N., 2004. Central limit theorem for sequential monte carlo methods
936 and its application to bayesian inference. *The Annals of Statistics* 32 (6),
937 pp. 2385–2411.
- 938 Chopin, N., Jacob, P. E., Papaspiliopoulos, O., 2013. Smc2: an efficient
939 algorithm for sequential analysis of state space models. *Journal of the*
940 *Royal Statistical Society: Series B (Statistical Methodology)* 75 (3), 397
941 – 426.
- 942 Craig, S., Lewandowsky, S., Little, D. R., May 2011. Error discounting
943 in probabilistic category learning. *Journal of experimental psychology.*
944 *Learning, memory, and cognition* 37 (3), 673–87.
- 945 Del Moral, P., 1996. Non-linear filtering: interacting particle resolution.
946 *Markov processes and related fields* 2 (4), 555–581.
- 947 Del Moral, P., 2013. Mean field simulation for Monte Carlo integration.
948 Chapman & Hall/CRC Press, London.
- 949 Douc, R., Cappé, O., Moulines, E., 2005. Comparison of resampling schemes
950 for particle filtering. In: *Proceedings of the 4th International Symposium*

951 on Image and Signal Processing and Analysis. ISPA 2005. IEEE, pp. 64–
952 69.

953 Douc, R., Garivier, A., Moulines, E., Olsson, J., 2011. Sequential Monte
954 Carlo smoothing for general state space hidden Markov models. *The An-*
955 *nals of Applied Probability* 21, 2109–2145.

956 Douc, R., Moulines, E., 10 2008. Limit theorems for weighted samples with
957 applications to sequential Monte Carlo methods. *Ann. Statist.* 36 (5),
958 2344–2376.

959 Doucet, A., de Freitas, N., Gordon, N., 2001a. An introduction to sequential
960 monte carlo methods. In: Doucet, A., de Freitas, N., Gordon, N. (Eds.),
961 *Sequential Monte Carlo methods in practice*. Springer, New York, pp.
962 3–14.

963 Doucet, A., de Freitas, N., Gordon, N. (Eds.), 2001b. *Sequential Monte*
964 *Carlo methods in practice*. Springer, New York.

965 Doucet, A., Johansen, A. M., 2011. A tutorial on particle filtering and
966 smoothing: Fifteen years later. In: Crisan, D., Rozovskii, B. (Eds.), *The*
967 *Oxford handbook of nonlinear filtering*. Oxford University Press, New
968 York, NY, pp. 656–704.

969 Gilks, W. R., Berzuini, C., 2001. Following a moving target: Monte Carlo
970 inference for dynamic Bayesian models. *Journal of the Royal Statistical*
971 *Society. Series B (Statistical Methodology)* 63 (1), pp. 127–146.

972 Gordon, N. J., Salmond, D. P., Smith, A. F. M., 1993. A novel approach to
973 non-linear/non-Gaussian Bayesian state estimation. *IEE Proceedings-F*
974 140, 107–113.

- 975 Isard, M., Blake, A., 1998. Condensation—conditional density propagation
976 for visual tracking. *International Journal of Computer Vision* 29 (1), 5–28.
977 URL <http://dx.doi.org/10.1023/A:1008078328650>
- 978 Johansen, A. M., Doucet, A., Davy, M., 2008. Particle methods for maxi-
979 mum likelihood estimation in latent variable models. *Statistics and Com-
980 puting* 18 (1), 47–57.
- 981 Kahn, H., Marshall, A. W., 1953. Methods of reducing sample size in monte
982 carlo computations. *Journal of the Operations Research Society of Amer-
983 ica* 1 (5), 263–278.
- 984 Kalman, R. E., 1960. A new approach to linear filtering and prediction
985 problems. *Transactions of the American Society of Mechanical Engineers,
986 Series D, Journal of Basic Engineering* 82, 35–45.
- 987 Kalman, R. E., Bucy, R. S., 1961. New results in linear filtering and pre-
988 diction theory. *Transactions of the American Society of Mechanical Engi-
989 neers, Series D, Journal of Basic Engineering* 83, 95–108.
- 990 Kantas, N., Doucet, A., Singh, S. S., Maciejowski, J. M., Chopin, N., 2015.
991 On particle methods for parameter estimation in general state-space mod-
992 els. *Statistical Science* 30 (3), 328–351.
- 993 Kitagawa, G., 1996. Monte carlo filter and smoother for non-gaussian non-
994 linear state space models. *Journal of Computational and Graphical Statis-
995 tics* 5, 1–25.
- 996 Kitagawa, G., Sato, S., 2001. Monte carlo smoothing and self-organizing
997 state?space model. In: Doucet, A., de Freitas, J. F. G., Gordon, N. J.

- 998 (Eds.), Sequential Monte Carlo methods in practice. Springer, New York,
999 pp. 177–195.
- 1000 Knowlton, B. J., Squire, L. R., Gluck, M. A., 1994. Probabilistic classifica-
1001 tion learning in amnesia. *Learning & Memory* 1, 106–120.
- 1002 Kruschke, J. K., 1992. ALCOVE: an exemplar-based connectionist model of
1003 category learning. *Psychological Review* 99, 22–44.
- 1004 Liu, J., West, M., 2001. Combined parameter and state estimation in
1005 simulation-based filtering. In: Doucet, A., de Freitas, N., Gordon, N.
1006 (Eds.), Sequential Monte Carlo Methods in Practice. Springer, New York,
1007 NY.
- 1008 Liu, J. S., 2001. Monte Carlo Strategies in Scientific Computing. Springer,
1009 New York.
- 1010 Liu, J. S., Chen, R., 1998. Sequential Monte Carlo methods for dynamic
1011 systems. *Journal of the American Statistical Association* 93 (443), 1032–
1012 1044.
- 1013 Montemerlo, M., Thrun, S., Koller, D., Wegbreit, B., 2002. Fastslam: A
1014 factored solution to the simultaneous localization and mapping problem.
1015 In: Dechter, R., Kearns, M., Sutton, R. (Eds.), Eighteenth National Con-
1016 ference on Artificial Intelligence. American Association for Artificial In-
1017 telligence, Menlo Park, CA, USA, pp. 593–598.
- 1018 Myung, J. I., Cavagnaro, D. R., Pitt, M. A., 2013. A tutorial on adaptive
1019 design optimization. *Journal of Mathematical Psychology* 57 (3–4), 53–67.
- 1020 Nummiaro, K., Koller-Meier, E., Gool, L. V., 2003. An adaptive color-based

- 1021 particle filter. *Image and Vision Computing* 21 (1), 99 – 110.
1022 URL <http://www.sciencedirect.com/science/article/pii/S0262885602001294>
- 1023 Pitt, M. K., Shephard, N., 1999. Filtering via simulation: Auxiliary particle
1024 filters. *Journal of the American Statistical Association* 94 (446), 590.
- 1025 Robert, C. P., Casella, G., 2004. *Monte Carlo statistical methods*, 2nd Edi-
1026 tion. Springer, New York.
- 1027 Sanborn, A. N., Griffiths, T. L., Navarro, D. J., 2010. Rational approxi-
1028 mations to rational models: Alternative algorithms for category learning.
1029 *Psychological Review* 117 (4), 1144–1167.
- 1030 Speekenbrink, M., Channon, S., Shanks, D. R., 2008. Learning strategies in
1031 amnesia. *Neuroscience and Biobehavioral Reviews* 32, 292–310.
- 1032 Speekenbrink, M., Shanks, D. R., May 2010. Learning in a changing envi-
1033 ronment. *Journal of Experimental Psychology. General* 139 (2), 266–98.
- 1034 Van Zandt, T., Sep. 2000. How to fit a response time distribution. *Psycho-*
1035 *nomic Bulletin & Review* 7 (3), 424–65.
- 1036 Visser, I., 2011. Seven things to remember about hidden Markov models:
1037 A tutorial on Markovian models for time series. *Journal of Mathematical*
1038 *Psychology* 55 (6), 403–415.
- 1039 Whiteley, N., 2013. Stability properties of some particle filters. *The Annals*
1040 *of Applied Probability* 23 (6), 2500–2537.
- 1041 Whiteley, N., Johansen, A. M., 2011. *Bayesian time series models*. Cam-
1042 bridge University Press, Cambridge, U.K., Ch. Auxiliary particle filtering:
1043 Recent developments, pp. 52–81.

1044 Yi, M. S., Steyvers, M., Lee, M., 2009. Modeling human performance in
1045 restless bandits with particle filters. *The Journal of Problem Solving* 2 (2),
1046 5.