
Inference Networks for Sequential Monte Carlo in Graphical Models

Brooks Paige
Frank Wood

Department of Engineering Science, University of Oxford

BROOKS@ROBOTS.OX.AC.UK
FWOOD@ROBOTS.OX.AC.UK

Abstract

We introduce a new approach for amortizing inference in directed graphical models by learning heuristic approximations to stochastic inverses, designed specifically for use as proposal distributions in sequential Monte Carlo methods. We describe a procedure for constructing and learning a structured neural network which represents an inverse factorization of the graphical model, resulting in a conditional density estimator that takes as input particular values of the observed random variables, and returns an approximation to the distribution of the latent variables. This recognition model can be learned offline, independent from any particular dataset, prior to performing inference. The output of these networks can be used as automatically-learned high-quality proposal distributions to accelerate sequential Monte Carlo across a diverse range of problem settings.

1. Introduction

Recently proposed methods for Bayesian inference based on sequential Monte Carlo (Doucet et al., 2001) have shown themselves to provide state-of-the-art results in applications far broader than the traditional use of sequential Monte Carlo (SMC) for filtering in state space models (Gordon et al., 1993; Pitt and Shephard, 1999), with diverse application to factor graphs (Naesseth et al., 2014), hierarchical Bayesian models (Lindsten et al., 2014), procedural generative graphics (Ritchie et al., 2015), and general probabilistic programs (Wood et al., 2014; Todeschini et al., 2014). These are accompanied by complementary computational advances, including memory-efficient implementations (Jun and Bouchard-Côté, 2014), and highly-parallel variants (Murray et al., 2014; Paige et al., 2014).

All these algorithms, however, share the need for specifying

a series of *proposal distributions*, used to sample candidate values at each stage of the algorithm. Sequential Monte Carlo methods perform inference progressively, iteratively targeting a sequence of intermediate distributions which culminates in a final target distribution. Well-chosen proposal distributions for transitioning from one intermediate target distribution to the next can lead to sample-efficient inference, and are necessary for practical application of these methods to difficult inference problems. Theoretically optimal proposal distributions (Doucet et al., 2000; Cornebise et al., 2008) are in general intractable, thus in practice implementing these algorithms requires either active (human) work to design an appropriate proposal distribution prior to sampling, or using an online estimation procedure to approximate the optimal proposal during inference (as in e.g. Van Der Merwe et al. (2000) or Cornebise et al. (2014) for state-space models). In many cases, a baseline proposal distribution which simulates from a prior distribution can be used, analogous to the so-called bootstrap particle filter for inference in state-space models; however, when confronted with tightly peaked likelihoods (i.e. highly informative observations), proposing from the prior distribution may be arbitrarily statistically inefficient (Del Moral and Murray, 2015). Furthermore, for some choices of sequences of densities there is no natural prior distribution, or even it may not be available in closed form. All in all, the need to design appropriate proposal distributions is a real impediment to the automatic application of these SMC methods to new models and problems.

This paper investigates how autoregressive neural network models for modeling probability distributions (Bengio and Bengio, 1999; Uria et al., 2013; Germain et al., 2015) can be leveraged to automate the design of model-specific proposal distributions for sequential Monte Carlo. We propose a method for learning proposal distributions for a given probabilistic generative model offline, prior to performing inference on any particular dataset. The learned proposals can then be reused as desired, allowing SMC inference to be performed quickly and efficiently for the same probabilistic model, but for new data — that is, for new settings of the observed random variables — once we have incurred the up-front cost of learning the proposals.

We thus present this work as an amortized inference procedure in the sense of Gershman and Goodman (2014), in that it takes a model as its input and generates an artifact which then can be leveraged for accelerating future inference tasks. Such procedures have been considered for other inference methods: learning idealized Gibbs samplers offline for models in which closed-form full conditionals are not available (Stuhmüller et al., 2013), using pre-trained neural networks to inform local MCMC proposal kernels (Jampani et al., 2015; Kulkarni et al., 2015), and learning messages for new factors for expectation-propagation (Heess et al., 2013). In the context of SMC, offline learning of high-quality proposal distributions provides a similar opportunity for amortizing runtime costs of inference, while simultaneously automating a currently-manual process.

2. Preliminaries

A directed graphical model, or Bayesian network (Pearl and Russell, 1998), defines a joint probability distribution and conditional independence structure via a directed acyclic graph. For each x_i in a set of random variables x_1, \dots, x_N , the network structure specifies a conditional density $p_i(x_i|\text{PA}(x_i))$, where $\text{PA}(x_i)$ denotes the parent nodes of x_i . Inference tasks in Bayesian networks involve marking certain nodes as observed random variables, and characterizing the posterior distribution of the remaining latent nodes. The joint distribution over N latent random variables \mathbf{x} and M observed random variables \mathbf{y} is defined as

$$p(\mathbf{x}, \mathbf{y}) \triangleq \prod_{i=1}^N f_i(x_i|\text{PA}(x_i)) \prod_{j=1}^M g_j(y_j|\text{PA}(y_j)), \quad (1)$$

where f_i and g_j refer to the probability density or mass functions associated with the respective latent and observed random variables x_i, y_j . Posterior inference in directed graphical models entails using Bayes' rule to estimate the posterior distribution of the latent variables \mathbf{x} given particular observed values \mathbf{y} ; that is, to characterize the target density $\pi(\mathbf{x}) \equiv p(\mathbf{x}|\mathbf{y})$. In most models, exact posterior inference is intractable, and one must resort to either variational or finite-sample approximations.

2.1. Sequential Monte Carlo

Importance sampling methods approximate expectations with respect to a (presumably intractable) distribution $\pi(\mathbf{x})$ by weighting samples drawn from a (presumably simpler) proposal distribution $q(\mathbf{x})$. In graphical models, with $\pi(\mathbf{x}) \equiv p(\mathbf{x}|\mathbf{y})$, we define an unnormalized target density $\gamma(\mathbf{x}) \equiv p(\mathbf{x}, \mathbf{y})$ such that $\pi(\mathbf{x}) = Z^{-1}\gamma(\mathbf{x})$, where the normalizing constant Z is unknown.

The sequential Monte Carlo algorithms we consider (Doucet

et al., 2001) for inference on an N -dimensional latent space $\mathbf{x}_{1:N}$ proceed by incrementally importance sampling a weighted set of K particles, with interspersed resampling steps to direct computation towards more promising regions of the high-dimensional space. We break the problem of estimating the posterior distribution of $\mathbf{x}_{1:N}$ into a series of simpler lower-dimensional problems by constructing an artificial sequence of target densities π_1, \dots, π_N (and corresponding unnormalized densities $\gamma_1, \dots, \gamma_N$) defined on increasing subsets $\mathbf{x}_{1:n}$, $n = 1, \dots, N$, where the final $\pi_N \equiv \pi$ is the full target posterior of interest. At each intermediate density, the importance sampling density $q_{n+1}(x_{n+1}|\mathbf{x}_{1:n})$ only needs to adequately approximate a low-dimensional step from $\mathbf{x}_{1:n}$ to x_{n+1} .

Procedurally, we initialize at $n = 1$ by sampling K values of x_1 from a proposal density $q_1(x_1)$, and assigning each of these particles x_1^k an associated importance weight

$$w_1(x_1^k) = \frac{\gamma_1(x_1^k)}{q_1(x_1^k)}, \quad W_1^k = \frac{w_1(x_1^k)}{\sum_{j=1}^K w_1(x_1^j)}. \quad (2)$$

For each subsequent $n = 2, \dots, N$, we first resample the particles according to the normalized weights at W_{n-1}^k , preferentially duplicating high-weight particles and discarding those with low weight. To do this we draw particle ancestor indices $a_{n-1}^1, \dots, a_{n-1}^K$ from a resampling distribution $r(\cdot|W_{n-1}^1, \dots, W_{n-1}^K)$ corresponding to any standard resampling scheme (Douc et al., 2005). We then extend each particle by sampling a value for x_n^k from the proposal kernel $q_n(x_n^k|\cdot)$, and update the importance weights

$$w_n(\mathbf{x}_{1:n}^k) = \frac{\gamma_n(\mathbf{x}_{1:n}^k)}{\gamma_{n-1}(\mathbf{x}_{1:n-1}^{a_{n-1}^k})q_n(x_n^k|\mathbf{x}_{1:n-1}^{a_{n-1}^k})}, \quad (3)$$

$$W_n^k = \frac{w_n(\mathbf{x}_{1:n}^k)}{\sum_{j=1}^K w_n(\mathbf{x}_{1:n}^j)}. \quad (4)$$

We can approximate expectations with respect to the target density $\pi(\mathbf{x}_{1:N})$ using the SMC estimator

$$\hat{\pi}(\mathbf{x}_{1:N}^{1:K}) = \sum_{k=1}^K W_N^k \delta_{\mathbf{x}_{1:N}^k}(\mathbf{x}_{1:N}), \quad (5)$$

where $\delta(\cdot)$ is a Dirac point mass.

2.2. Target densities and proposal kernels

The choice of incremental target densities is application-specific; innovation in SMC algorithms often involves proposing novel manners for constructing sequences of intermediate distributions. These incremental densities do not necessarily need to correspond to marginal distributions of full target. Particularly relevant recent work directed towards improving SMC inference in the same class of models we address includes the Biips ordering and arrangement

algorithm (Todeschini et al., 2014), the divide and conquer approach (Lindsten et al., 2014), and heuristics for scoring orderings in general factor graphs (Naesseth et al., 2014; 2015). All these methods provide a means for selecting a sequence of intermediate target densities — however, given a sequence of targets, one still must supply an appropriate proposal density.

The ideal choice for this proposal in general is found by proposing directly from the incremental change in densities (Doucet et al., 2000), with

$$q_n^{opt}(x_n|x_{1:n-1}) = \frac{\pi_n(x_{1:n})}{\pi_{n-1}(x_{1:n-1})} \propto \frac{\gamma_n(x_{1:n})}{\gamma_{n-1}(x_{1:n-1})}. \quad (6)$$

Using this proposal, each of the unnormalized weights in Equation (3) are independent of the sampled values of x_n^k . In practice this conditional density is nearly always intractable, and one must resort to approximation.

Adaptive importance sampling methods aim to learn the optimal proposal online during the course of inference, immediately prior to proposing values for the next target density. In both in the context of population Monte Carlo (PMC) (Cappé et al., 2008) and sequential Monte Carlo (Cornebise et al., 2008; 2014; Gu et al., 2015), a parametric family $q(\mathbf{x}|\lambda)$ is proposed, with λ is a free parameter, and the adaptive algorithms aim to minimize either the reverse Kullback-Leibler (KL) divergence or Chi-squared distance between the approximating family $q(\mathbf{x}|\lambda)$ and the optimal proposal density. This can be optimized via stochastic gradient descent (Gu et al., 2015), or for specific forms of q by online Monte Carlo expectation maximization, both for population Monte Carlo (Cappé et al., 2008) and in state-space models (Cornebise et al., 2014). Note that this is the reverse of the KL divergence traditionally used in variational inference (Jordan et al., 1999), and takes the form of an expectation with respect to the intractable target distribution.

2.3. Neural autoregressive distribution estimation

As a general model class for $q(\mathbf{x}|\cdot)$, we adapt recent advances in flexible neural network density estimators, appropriate for both discrete and continuous high-dimensional data. We focus particularly on the use of autoregressive neural network density estimation models (Bengio and Bengio, 1999; Larochelle and Murray, 2011; Uria et al., 2013; Germain et al., 2015) which model high-dimensional distribution by learning a sequence of one-dimensional conditional distributions; that is, learning each product term in

$$p(\mathbf{x}) = \prod_{n=1}^N p(x_n|x_1, \dots, x_{n-1}), \quad (7)$$

typically with weight parameter sharing across densities.

We choose to adapt the masked autoencoder for distribution estimation (MADE) model (Germain et al., 2015), which

fits an autoregressive model to binary data, with structure inspired by autoencoders. In its simplest form, a single-layer MADE model described on N -dimensional binary data $\mathbf{x} \in [0, 1]^N$ has a hidden layer $\mathbf{h}(\mathbf{x})$ and output $\hat{\mathbf{x}}$ with

$$\mathbf{h}(\mathbf{x}) = \sigma_w(\mathbf{b} + (\mathbf{W} \odot \mathbf{M}_w)\mathbf{x}) \quad (8)$$

$$\hat{\mathbf{x}} = \sigma_v(\mathbf{c} + (\mathbf{V} \odot \mathbf{M}_v)\mathbf{h}(\mathbf{x})), \quad (9)$$

where \mathbf{b} , \mathbf{c} , \mathbf{W} , \mathbf{V} are real-valued parameters to be learned, \odot denotes elementwise multiplication, σ_w, σ_v are nonlinear functions, and $\mathbf{M}_w, \mathbf{M}_v$ are fixed binary masks. Critically, the construction of the masks is such that computing the network output for each \hat{x}_n requires only the inputs x_1, \dots, x_{n-1} , with the zeros in the masks dropping the connections. The masks are generated by assigning each unit in each hidden layer a number from $1, \dots, N-1$, describing which of the dimensions x_1, \dots, x_{n-1} it is permitted to take as input; output units then are only permitted to take as input hidden nodes numbered lower than their output.

With a logistic function sigmoid as σ_v , then \hat{x}_n can be interpreted as $p(x_n|x_1, \dots, x_{n-1})$, and to compute \hat{x}_n one does not need supply any value as input to $\mathbf{h}(\mathbf{x})$ for the dimensions x_n, \dots, x_N . That is, if one follows all connections “back” through the network from \hat{x}_n to the input \mathbf{x} , one would find only themselves at x_1, \dots, x_{n-1} .

3. Approach

Our approach is two-fold. First, given a Bayesian network that acts as a generative model for our observed data \mathbf{y} given latent variables \mathbf{x} , we construct a new Bayesian network which acts as a generative model for our latent \mathbf{x} , given observed data \mathbf{y} . This network is constructed such that the joint distribution of the new “inverse model”, which we will refer to as $\tilde{p}(\mathbf{x}, \mathbf{y}) = \tilde{p}(\mathbf{y})\tilde{p}(\mathbf{x}|\mathbf{y})$, preserves the conditional dependence structure in the original model $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$, but has a different factorization (Stuhlmüller et al., 2013).

Unfortunately, unlike the original forward model, the inverse model has conditional densities which we do not in general know how to normalize or sample from. However, were we to know the full conditional density of the inverse model $\tilde{p}(\mathbf{x}|\mathbf{y})$, then we could directly draw samples of \mathbf{x} given a particular dataset \mathbf{y} .

Thus second task is to learn approximations for the conditionals $\tilde{p}(x_i|\tilde{\text{PA}}(x_i))$, where $\tilde{\text{PA}}(x_i)$ are parents of x_i in the inverse model. To do so we employ neural density estimators and design a procedure to train these “offline”, in the sense that no real data is required.

3.1. Defining the inverse model

We begin by constructing an inverse model $\tilde{p}(\mathbf{x}, \mathbf{y})$ which admits the same distribution over all random variables as

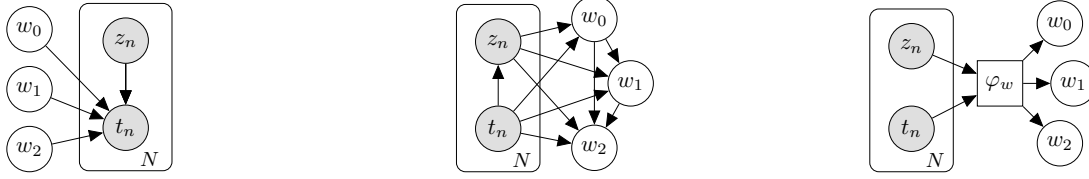


Figure 1. A non-conjugate regression model, as (left) a Bayes net representing a generative model for the data $\{t_n\}$; (middle) with dependency structure inverted, as a generative model for the latent variables w_0, w_1, w_2 ; (right) showing the explicit neural network structure of the learned approximation to the inverse conditional distribution $\tilde{p}(w_{0:2}|z_{1:N}, t_{1:N})$. New datasets $\{z_n, t_n\}_{n=1}^N$ can be input directly into the joint density estimator φ_w to estimate the posterior. Note that the ordering of the latent variables $w_{0:2}$ used in this example is chosen arbitrarily; any permutation of the latent variables would not change the overall structure of the inverse model.

$p(\mathbf{x}, \mathbf{y})$, but with a different factorization. We first note that the directed acyclic graph structure of $p(\mathbf{x}, \mathbf{y})$ imposes a partial ordering on all random variables \mathbf{x} and \mathbf{y} ; we choose any single valid ordering arbitrarily, and define the sequences x_1, \dots, x_N and y_1, \dots, y_M such that for any x_i , $\text{PA}(x_i) \subseteq \{x_1, \dots, x_{i-1}\} \cup \{y_j\}_{j=1}^M$, and for any y_j , $\text{PA}(y_j) \subseteq \{y_1, \dots, y_{j-1}\} \cup \{x_i\}_{i=1}^N$.

Our goal here is to construct as simple as possible a distribution $\tilde{p}(\mathbf{x}|\mathbf{y})$ whose factorization does not introduce any new conditional independencies not also present in the original generative model. To consider two extremes: a fully factorized $\tilde{p}(\mathbf{x}|\mathbf{y}) \equiv \prod_{i=1}^N \tilde{p}(x_i|\mathbf{y})$ which assumes all x_i are conditionally independent given \mathbf{y} may be attractive for computational reasons, but fails to capture all the structure of the posterior; whereas a fully connected $\tilde{p}(\mathbf{x}|\mathbf{y}) \equiv \prod_{i=1}^N \tilde{p}(x_i|x_{1:i-1}, \mathbf{y})$ is guaranteed to capture all dependencies, but may be unnecessarily complex.

To define the approximating distribution at each x_i , we invert the dependencies on y_j , effectively running the generative model backwards. Following the heuristic algorithm of Stuhlmüller et al. (2013), we do this by literally constructing the dependency graph in reverse. Ordering the random variables $y_M, \dots, y_1, x_N, \dots, x_1$, we define a new parent set $\tilde{\text{PA}}(x_i)$ for each x_i in the transformed model, with $\tilde{\text{PA}}(x_i) \subseteq \{x_{i+1}, \dots, x_N, y_1, \dots, y_M\}$. In particular, if $\text{MB}(x_i)$ is the Markov blanket of the latent variable x_i in the original probability model, then defining the parent sets

$$\begin{aligned} \tilde{\text{PA}}(x_i) &= \text{MB}(x_i) \cap \{x_{i+1}, \dots, x_N, y_1, \dots, y_M\} \\ \tilde{\text{PA}}(y_j) &= \text{MB}(y_j) \cap \{y_{j+1}, \dots, y_M\} \end{aligned}$$

yields a model with the same dependency structure as the original model $p(\mathbf{x}, \mathbf{y})$; however, now the sequence is reversed such that the observed values are inputs (i.e., $\tilde{\text{PA}}(y_j) \cap \mathbf{x} = \emptyset$). The sequence under the new model, which we will refer to as $\tilde{p}(\mathbf{x}, \mathbf{y})$, factorizes naturally as $\tilde{p}(\mathbf{x}, \mathbf{y}) = \tilde{p}(\mathbf{x}|\mathbf{y})\tilde{p}(\mathbf{y})$; particularly important to us is the factorization of the conditional density $\tilde{p}(\mathbf{x}|\mathbf{y}) = \prod_{i=1}^N \tilde{p}(x_i|\tilde{\text{PA}}(x_i))$.

This algorithm produces inverse graph structures which despite not being fully connected, are guaranteed to represent all conditional dependencies in the original graph:

Proposition 1. Preservation of conditional dependence. Let x_A, x_B, x_C be any latent or observed random variables in $p(\mathbf{x})$ with graph structure G , and let $\tilde{x}_A, \tilde{x}_B, \tilde{x}_C$ denote the corresponding random variables in the inverse model $\tilde{p}(\mathbf{x})$ with graph structure \tilde{G} , constructed via the algorithm above. Then if \tilde{x}_A and \tilde{x}_B are conditionally independent given \tilde{x}_C in the inverse model \tilde{G} , they were also conditionally independent in the original model G ; that is,

$$\tilde{x}_A \perp\!\!\!\perp \tilde{x}_B | \tilde{x}_C \Rightarrow x_A \perp\!\!\!\perp x_B | x_C.$$

Proof. Suppose we had a conditional dependence in G which was not preserved in \tilde{G} , i.e. with $x_A \not\perp\!\!\!\perp x_B | x_C$ but $\tilde{x}_A \perp\!\!\!\perp \tilde{x}_B | \tilde{x}_C$. Without loss of generality assume \tilde{x}_B was added to the inverse graph prior to \tilde{x}_A , i.e. $x_A \prec x_B$ in G . Note that $x_A \not\perp\!\!\!\perp x_B | x_C$ can occur either due to a direct dependence between x_A and x_B , or, due to both $x_A, x_B \in \text{PA}(x_C)$; in either case, $x_B \in \text{MB}(x_A)$. Then when adding \tilde{x}_B to the inverse graph \tilde{G} we are guaranteed to have $\tilde{x}_B \in \tilde{\text{PA}}(\tilde{x}_A)$, in which case $\tilde{x}_A \not\perp\!\!\!\perp \tilde{x}_B$. \square

Examples of generative models and their corresponding inverse models are shown in Figures 1–3. Note that as the topological sort of the nodes in the original generative model is not unique, neither is the inverse graphical model.

3.2. Learning a family of approximating densities

Following Cappé et al. (2008), learning proposals for importance sampling on $\pi(\mathbf{x})$ in a single-dataset setting (i.e., with fixed \mathbf{y}) entails proposing a parametric family $q(\mathbf{x}|\lambda)$, where λ is a free parameter, and then choosing λ to minimize

$$D_{KL}(\pi||q_\lambda) = \int \pi(\mathbf{x}) \log \left[\frac{\pi(\mathbf{x})}{q(\mathbf{x}|\lambda)} \right] d\mathbf{x}. \quad (10)$$

This KL divergence between the true posterior distribution $\pi(\mathbf{x}) \equiv p(\mathbf{x}|\mathbf{y})$ and proposal distribution $q(\mathbf{x}|\lambda)$ is also known as the relative entropy criterion, and is a preferred objective function in situations in which the estimation goal construct a high-quality weighted sample representation, rather than to minimize the variance of a particular expectation (Cornebise et al., 2008).

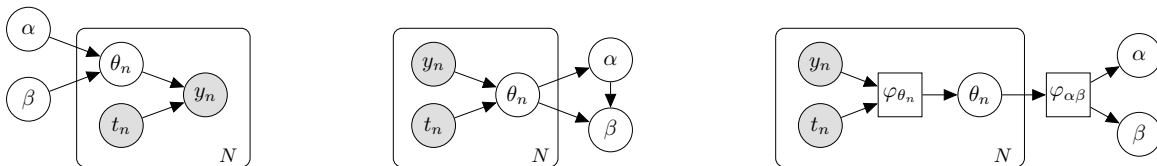


Figure 2. A hierarchical Bayesian model. (left) A generative model for the data $\{x_n\}$; (middle) with dependency structure inverted; (right) showing the two distinct joint neural conditional density estimators. Note in particular the inverse model still partially factorizes across the latent variables. The learned factor φ_{θ_n} is replicated N times in the inverse model, allowing re-use of weights, simplifying training.

In an amortized inference setting, instead of learning λ explicitly for a fixed value of \mathbf{y} , we learn a mapping from \mathbf{y} to λ . More explicitly, if $\mathbf{y} \in \mathcal{Y}$ and $\lambda \in \vartheta$, then learning a deterministic mapping $\varphi : \mathcal{Y} \rightarrow \vartheta$ allows performing approximate inference for $p(\mathbf{x}|\mathbf{y})$ with only the computational complexity of evaluating the function φ . The tradeoff is that the training of φ itself may be quite involved.

We thus generalize the adaptive importance sampling algorithms by learning a family of distributions $q(\mathbf{x}|\mathbf{y})$, parameterized by the observed data \mathbf{y} . Suppose that $\lambda = \varphi(\eta, \mathbf{y})$, where the function φ is parameterized by a set of upper-level parameters η . We would like a choice of η which performs well across all datasets \mathbf{y} . We can frame this as minimizing the expected value of Eq. (10) under $p(\mathbf{y})$, suggesting an objective function $\mathcal{J}(\eta)$ defined as

$$\begin{aligned} \mathcal{J}(\eta) &= \int D_{KL}(\pi||q_\lambda)p(\mathbf{y})d\mathbf{y} \\ &= \int p(\mathbf{y}) \int p(\mathbf{x}|\mathbf{y}) \log \left[\frac{p(\mathbf{x}|\mathbf{y})}{q(\mathbf{x}|\varphi(\eta, \mathbf{y}))} \right] d\mathbf{x}d\mathbf{y} \\ &= \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [-\log q(\mathbf{x}|\varphi(\eta, \mathbf{y}))] + \text{const} \end{aligned} \quad (11)$$

which has a gradient

$$\nabla_\eta \mathcal{J}(\eta) = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [-\nabla_\eta \log q(\mathbf{x}|\varphi(\eta, \mathbf{y}))]. \quad (12)$$

Notice that these expectations in Equations (11) and (12) are with respect to the tractable joint distribution $p(\mathbf{x}, \mathbf{y})$. We can thus fit η by stochastic gradient descent, estimating the expectation of the gradient $\nabla_\eta \mathcal{J}(\eta)$ by sampling synthetic full-data training examples $\{\mathbf{x}, \mathbf{y}\}$ from the original model. This procedure can be performed entirely offline — we require only to be able to sample from the joint distribution $p(\mathbf{x}, \mathbf{y})$ to generate candidate data points (effectively providing infinite training data). In any directed graphical model this can be achieved by ancestral sampling, where in addition to sampling \mathbf{x} we sample values of the as-yet unobserved variables \mathbf{y} . Furthermore, we do not need need to be able to compute gradients of our model $p(\mathbf{x}, \mathbf{y})$ itself — we only need the gradients of our recognition model $q(\mathbf{x}|\varphi(\eta, \mathbf{y}))$, allowing use of any differentiable representation for q . We choose the parametric family $q(\mathbf{x}|\lambda)$ and the transformation φ such that this inner gradient in Eq. (12) can be computed easily.

We can now use the conditional independence structure in our inverse model $\tilde{p}(\mathbf{x}, \mathbf{y})$ to break down $q(\mathbf{x}|\lambda)$, an approximation of $\tilde{p}(\mathbf{x}|\mathbf{y})$, into a product of smaller conditional densities each approximating $\tilde{p}(x_i|\widetilde{\text{PA}}(x_i))$. The full proposal density $q(\mathbf{x}|\varphi(\eta, \mathbf{y}))$ can be decomposed as

$$q(\mathbf{x}|\varphi(\eta, \mathbf{y})) = \prod_{i=1}^N q_i(x_i|\varphi_i(\eta_i, \widetilde{\text{PA}}(x_i))) \quad (13)$$

with the gradient similarly decomposing as

$$\nabla_{\eta_i} \mathcal{J}(\eta) = \mathbb{E}_{p(\mathbf{x}, \mathbf{y})} [-\nabla_{\eta_i} \log q_i(x_i|\varphi_i(\eta_i, \widetilde{\text{PA}}(x_i)))].$$

Each of these expectations requires only samples of the random variables in $\{x_i\} \cup \widetilde{\text{PA}}(x_i)$, reducing the dimensionality of the joint optimization problem. This factorization of $q(\mathbf{x}|\varphi(\eta, \mathbf{y}))$ does not practically reduce the expressivity of the approximating family, as all conditional dependencies in the true posterior are preserved.

3.3. Joint conditional neural density estimation

We particularly wish to construct the inverse factorization $\tilde{p}(\mathbf{x}|\mathbf{y})$ (and our proposal model $q(\cdot)$) in such a way that we deal naturally with the presence of head-to-head nodes, in which one random variable may have a very large parent set. This situation is common in machine learning models: it is quite common to have generative models which factorize in the joint distribution, but have complex dependencies in the posterior; see for example the model in Figure 1.

We thus choose to treat all such situations in our inverse factorization — where a sequence of variables $\mathbf{x}' \subseteq \mathbf{x}$ are fully dependent on one another after conditioning on a shared set of parent nodes $\widetilde{\text{PA}}(\mathbf{x}')$ — as a single joint conditional density which we will approximate with an autoregressive density model. We extend MADE (Germain et al., 2015) to function as a conditional density estimator by allowing it to take $\widetilde{\text{PA}}(\mathbf{x}')$ as additional inputs, and constructing the masks such that these additional inputs are propagated through all hidden layers to all outputs, even for the very first dimension. As in MADE this can be achieved by labeling the hidden units with integers denoting which input dimensions they are allowed to accept. In contrast to the original MADE, we label hidden units with numbers from $0, \dots, N-1$, where

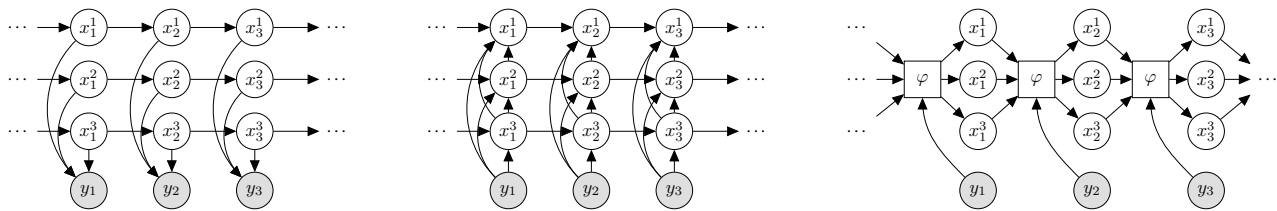


Figure 3. Factorial HMM. (left) The generative model consists D independent Markov models, with observed data y_t depending on the current state of each latent HMM. (middle) An inverse model obtained by reversing the order of the generative model at each t . Conditioned on the previous latent states at $t - 1$ and the next observation y_t , all latent states at each t are dependent on one another and must be modeled jointly. (right) The repeated structure at each $t = 1, 2, \dots$ means that the same learned conditional density network can be reused at every t .

hidden units labeled 0 to take as input only the dimensions in $\widetilde{\text{PA}}(\mathbf{x}')$. For single-dimensional data, where $N = 1$, all hidden units are labeled 0 and all feed forward into the single output x_1 , recovering a standard mixture density network (Bishop, 1994).

To model non-binary data, MADE can be extended by altering the output layer network to emit parameters of any univariate probability density function. We take the same approach by which RNADE (Uria et al., 2013) modifies the binary autoregressive distribution estimator NADE (Larochelle and Murray, 2011) to handle real-valued data, with an output layer that parameterizes a univariate mixture of D Gaussians for each dimension x_i conditioned on its parents. The probability of any particular x_i is given by

$$q(x_i | \varphi_i(\eta_i, \widetilde{\text{PA}}(x_i))) = \sum_{d=1}^D \alpha_{i,d} \mathcal{N}(x_i | \mu_{i,d}, \sigma_{i,d}^2)$$

where $\mathcal{N}(\cdot)$ is the Gaussian probability density. This requires an output layer with $3 \times D$ dimensions, to predict D each of means $\mu_{i,d}$, standard deviations $\sigma_{i,d}$, weights $\alpha_{i,d}$; to enforce positivity of standard deviations we apply a softplus function to the raw network outputs, and a softmax function to ensure $\alpha_{i,\cdot}$ is a probability vector.

3.4. Training the neural network

Contrary to many standard settings in which one is limited by the amount of data present, we are armed with a sampler $p(\mathbf{x}, \mathbf{y})$ which allows us to generate effectively infinite training data. This could be used to sample a ‘‘giant’’ synthetic dataset, which we then use for mini-batch training via gradient descent; however, then we must decide how large a dataset is required. Alternatively, we could sample a brand new set of training examples for every mini-batch, never re-using previous samples.

In testing we found that a hybrid training procedure, which samples new synthetic datasets based on performance on a held-out set of synthetic validation data, appeared more efficient than resampling a new synthetic dataset for each

new gradient update. We perform mini-batch gradient updates on η using synthetic training data, while evaluating on the validation set. If the validation error increases, or after a set maximum number of steps, we draw new sets of both synthetic training and validation data from $p(\mathbf{x}, \mathbf{y})$.

In all experiments we use Adam (Kingma and Ba, 2015) with the suggested default parameters to update learning rates online, and use rectified linear activation functions.

4. Examples

4.1. Inverting a single factor

To illustrate the basic method for inverting factors, we consider a non-conjugate polynomial regression model, with global-only latent variables. The graphical model, its inversion, and the neural network structure are shown in Figure 1. Here we place a Laplace prior on the regression weights, and have Student-t likelihoods, giving us

$$w_d \sim \text{Laplace}(0, 10^{1-d}) \quad \text{for } d = 0, 1, 2;$$

$$t_n \sim t_\nu(w_0 + w_1 z_n + w_2 z_n^2, \epsilon^2) \quad \text{for } n = 1, \dots, N$$

for fixed $\nu = 4, \epsilon = 1$, and $z_n \in (-10, 10)$ uniformly. The goal is to estimate the posterior distribution of weights for the constant, linear, and quadratic terms, given any possible collected dataset $\{z_n, t_n\}_{n=1}^N$. In the notation of the preceding sections, we have latent variables $\mathbf{x} \equiv \{w_0, w_1, w_2\}$ and observed variables $\mathbf{y} \equiv \{z_n, t_n\}_{n=1}^N$.

Note particularly that although the original graphical model which expressed $p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$ factorizes into products over y_n which are conditionally independent given \mathbf{x} , in the inverse model $\tilde{p}(\mathbf{x}|\mathbf{y})$ due to the explaining-away phenomenon all latent variables depend on all others: there are no latent variables which can be d -separated from the observed \mathbf{y} , and all latent variables share \mathbf{y} as parents. This means we fit as proposal only a single joint density $q(w_{0:2}|z_{1:N}, t_{1:N})$. Examples of representative output from this network are shown in Figure 4. The trained network used here 200 hidden units in each of two hidden layers, and a mixture of 3 Gaussians as each output.

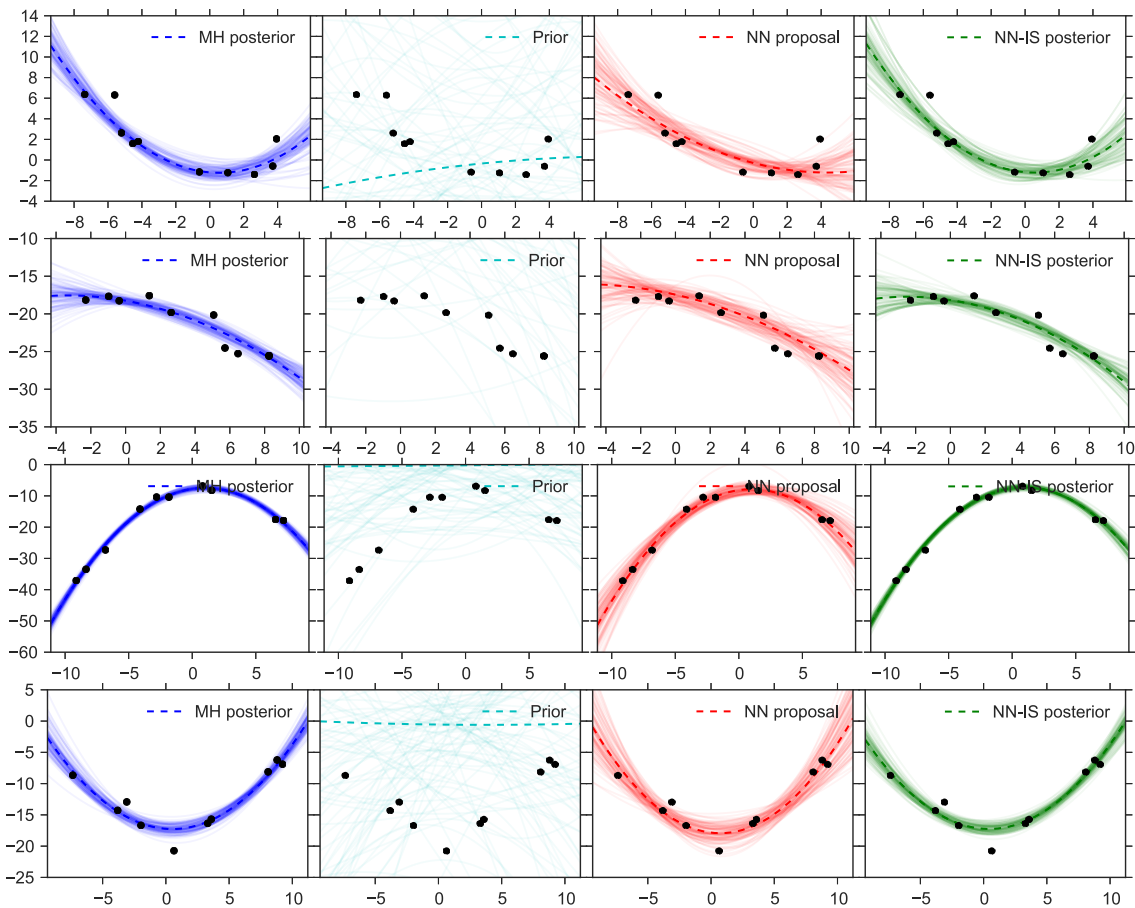


Figure 4. Representative output in the polynomial regression example. Plots show 100 samples each at 5% opacity, with the mean marked as a solid dashed line. These are all proposed using the same pre-trained neural network — not just the same neural network structure, but also identical learned weights. The MCMC posterior is generated by thinning 10000 samples by a factor 100, after 10000 samples of burnin. The neural network proposal yields estimated polynomial curves close to the true posterior solution, albeit slightly more diffuse.

4.2. A hierarchical Bayesian model

Consider as a new example a representative multilevel model where exact inference is intractable, a Poisson model for estimating failure rates of power plant pumps (George et al., 1993). Given N power plant pumps, each having operated for t_n thousands of hours, we see x_n failures, following

$$\begin{aligned} \alpha &\sim \text{Exponential}(1.0), & \beta &\sim \text{Gamma}(0.1, 1.0), \\ \theta_n &\sim \text{Gamma}(\alpha, \beta), & y_n &\sim \text{Poisson}(\theta_n t_n). \end{aligned}$$

The graphical model, an inverse factorization, and the neural network structure are shown in Figure 2. To generating synthetic training data, t_n are sampled *iid* from an exponential distribution with mean 50.

The repeated structure in the inverse factorization of this model allows us to learn a single inverse factor to represent the distribution $\tilde{p}(\theta_n | t_n, y_n)$ across all n . This yields a far simpler learning problem than were we forced to fit all of $\tilde{p}(\theta_{1:N} | t_{1:N}, y_{1:N})$ jointly. Further, the repeated structure

allows us to use a divide-and-conquer SMC algorithm (Lindsten et al., 2014) which works particularly efficiently on this model. Each of the N replicated structures are sampled in parallel with independent particle sets, weighted locally, and resampled; once all θ_n are sampled, we end by sampling α and β jointly, which need both be included in order to evaluate the final terms in the joint target density. We stress that there is no obvious baseline proposal density to use for a divide-and-conquer SMC algorithm, as neither the marginal prior nor posterior distributions over θ_n are available in closed form. Any usage of this algorithm requires manual specification of some proposal $q(\theta_n)$.

We test our proposals on the actual power pump failure data analyzed in George et al. (1993). The relative convergence speeds of marginal likelihood estimators from importance sampling from prior and neural network proposals, and SMC with neural network proposals, are shown in Figure 5. To capture the wide tails of the broad gamma distributions, we use a mixture of 10 Gaussians here at each output node, and 500 hidden units in each of two hidden layers.

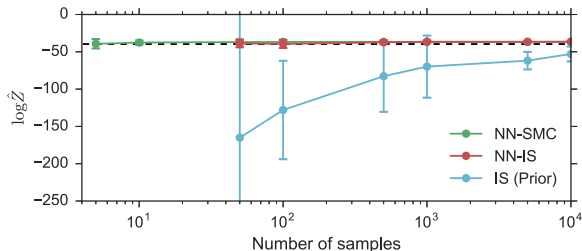


Figure 5. Convergence of marginal likelihood estimate as a function of number of particles, for likelihood-weighted importance sampling, neural network importance sampling, and a divide-and-conquer sequential Monte Carlo algorithm with neural network proposals. The SMC algorithm can achieve reasonable estimates of the normalizing constant with as few as 5 samples. Plot shows mean of 10 runs; error bars show two standard deviations.

4.3. Factorial hidden Markov model

Proposals can also be learned to approximate the optimal filtering distribution in models for sequential data; we demonstrate here on a factorial hidden Markov model (Ghahramani and Jordan, 1997), where each time step has a combinatorial latent space. The additive model we consider is inspired by the model studied in Kolter and Jaakkola (2012) for disaggregation of household energy usage; effective inference in this model is a subject of continued research. Some number of devices D are either in an active state, in which case each device i consumes μ^i units of energy, or it is off, in which case it consumes no energy. At each time step we receive a noisy observation of the total amount of energy consumed, summed across all devices. This model, whose graphical model structure is shown in Figure 3, can be represented as

$$\begin{aligned} x_t^i | x_{t-1}^i &\sim \text{Bernoulli}(\theta^i [x_{t-1}^i]) \\ y_t | x_t^1, \dots, x_t^D &\sim \mathcal{N}(\sum_{i=1}^D \mu^i x_t^i, \sigma^2), \end{aligned}$$

where θ^i represents the prior probability of devices switching on or off at each time increment. We design a synthetic example with $D = 20$, meaning each time step has $2^{20} \approx 100,000$ possible discrete states; the parameters μ^d are spread out from 30 to 500, with $\sigma = 10$. Each individual device has an initial probability 0.1 of being activated at $t = 1$, switching state at subsequent t with probability 0.05.

As different combinations of devices can yield identical total energy usage it is impossible to disambiguate between different combinations of active devices from a single observation, meaning any successful inference algorithm must attempt to mix across many disconnected modes over time to preserve the multiple possible explanations. Synthetic data and example output of inference is shown in the supplemental material. The effect of the learned proposals on the overall number of surviving particles is shown in Figure 6. Our proposal model uses D Bernoulli outputs in a 4-layer network, with 300 units per hidden layer; it takes as input

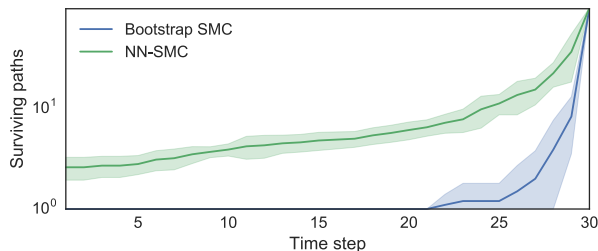


Figure 6. Learned proposals reduce particle degeneracy in the factorial HMM. Here we show the number of unique ancestries which survive over the course of 30 time steps, running 100 particles. Proposing from the transition dynamics nearly immediately degenerates to a single possible solution; the learned proposals increase the effective sample size at each stage and reduce the need for resampling. Plot shows mean and standard deviation over 10 runs.

the D latent states at the previous time $t - 1$, as well as the current observation y_t . A separate network is used for predicting the initial state x_1^i given only the initial input y_1 .

5. Discussion

We present this work primarily as a manner by which we compile away application-time inference costs when performing SMC, and automating the manual task of designing proposal densities. However, in some situations direct sampling from the model may provide a satisfactory approximation even eschewing importance weighting steps; in such cases our approach can be viewed as a graphical-model-regularized algorithm for designing and training neural networks with interpretable structural representations. Rather than learning from data, the emulator model is chosen to approximate the specified generative model, akin to the “sleep” cycle of the wake-sleep algorithm (Hinton et al., 1995).

In contrast to variational autoencoders (Kingma and Welling, 2014), where one simultaneously learns parameters for both the inference network and generative model from data, we assume a known generative model with fixed parameters and structured, interpretable latent variables. This provides robustness to bias arising from training data which comes from an unrepresentative sample, and also allows us to apply our method in situations where a sufficiently large supply of exemplar data is unavailable. However, it does require placing trust in the generative model: in particular, it requires a generative model which could plausibly create the data we will later collect and condition on.

Beyond these differences, our choice of $D_{KL}(\pi||q)$, the same minimized by EP, leads to approximations more appropriate for SMC refinement than a variational Bayes objective function; see e.g. Minka (2005) for a discussion of “zero-forcing” behavior, and e.g. Cappé et al. (2008) for a discussion of pathological cases in learned importance sampling distributions.

Acknowledgements

BP would like to thank both Tom Jin and Jan-Willem van de Meent for their helpful discussions, feedback, and ongoing commiseration. FW is supported under DARPA PPAML through the U.S. AFRL under Cooperative Agreement number FA8750-14-2-0006, Sub Award number 61160290-111668.

References

- Bengio, Y. and Bengio, S. (1999). Modeling high-dimensional discrete data with multi-layer neural networks. In *Advances in Neural Information Processing Systems*, volume 99, pages 400–406.
- Bishop, C. M. (1994). Mixture density networks. Technical report.
- Cappé, O., Douc, R., Guillin, A., Marin, J.-M., and Robert, C. P. (2008). Adaptive importance sampling in general mixture classes. *Statistics and Computing*, 18(4):447–459.
- Cornebise, J., Moulines, É., and Olsson, J. (2008). Adaptive methods for sequential importance sampling with application to state space models. *Statistics and Computing*, 18:461–480.
- Cornebise, J., Moulines, É., and Olsson, J. (2014). Adaptive sequential Monte Carlo by means of mixture of experts. *Statistics and Computing*, 24:317–337.
- Del Moral, P. and Murray, L. M. (2015). Sequential Monte Carlo with highly informative observations. *SIAM/ASA Journal on Uncertainty Quantification*, 3(1):969–997.
- Douc, R., Cappé, O., and Moulines, E. (2005). Comparison of resampling schemes for particle filtering. In *In 4th International Symposium on Image and Signal Processing and Analysis (ISPA)*, pages 64–69.
- Doucet, A., De Freitas, N., Gordon, N., et al. (2001). *Sequential Monte Carlo methods in practice*. Springer New York.
- Doucet, A., Godsill, S., and Andrieu, C. (2000). On sequential Monte Carlo sampling methods for Bayesian filtering. *Statistics and computing*, 10(3):197–208.
- George, E. I., Makov, U., and Smith, A. (1993). Conjugate likelihood distributions. *Scandinavian Journal of Statistics*, pages 147–156.
- Germain, M., Gregor, K., Murray, I., and Larochelle, H. (2015). MADE: masked autoencoder for distribution estimation. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015*, pages 881–889.
- Gershman, S. J. and Goodman, N. D. (2014). Amortized inference in probabilistic reasoning. In *Proceedings of the Thirty-Sixth Annual Conference of the Cognitive Science Society*.
- Ghahramani, Z. and Jordan, M. I. (1997). Factorial hidden Markov models. *Machine learning*, 29(2-3):245–273.
- Gordon, N. J., Salmond, D. J., and Smith, A. F. (1993). Novel approach to nonlinear/non-Gaussian Bayesian state estimation. *IEE Proceedings F (Radar and Signal Processing)*, 140(2):107–113.
- Gu, S., Ghahramani, Z., and Turner, R. E. (2015). Neural adaptive sequential Monte Carlo. In *Advances in Neural Information Processing Systems* 28.
- Heess, N., Tarlow, D., and Winn, J. (2013). Learning to pass expectation propagation messages. In *Advances in Neural Information Processing Systems* 26, pages 3219–3227.
- Hinton, G. E., Dayan, P., Frey, B. J., and Neal, R. M. (1995). The “wake-sleep” algorithm for unsupervised neural networks. *Science*, 268(5214):1158–1161.
- Jampani, V., Nowozin, S., Loper, M., and Gehler, P. V. (2015). The informed sampler: A discriminative approach to Bayesian inference in generative computer vision models. *Computer Vision and Image Understanding*, 136:32–44.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., and Saul, L. K. (1999). An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233.
- Jun, S.-H. and Bouchard-Côté, A. (2014). Memory (and time) efficient sequential Monte Carlo. In *Proceedings of the 31st international conference on Machine learning*, pages 514–522.
- Kingma, D. and Ba, J. (2015). Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Kingma, D. P. and Welling, M. (2014). Auto-encoding variational Bayes. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Kolter, J. Z. and Jaakkola, T. (2012). Approximate inference in additive factorial HMMs with application to energy disaggregation. In *International conference on artificial intelligence and statistics*, pages 1472–1482.
- Kulkarni, T. D., Kohli, P., Tenenbaum, J. B., and Mansinghka, V. K. (2015). Picture: a probabilistic programming language for scene perception. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.

- Larochelle, H. and Murray, I. (2011). The neural autoregressive distribution estimator. In *International Conference on Artificial Intelligence and Statistics*, pages 29–37.
- Lindsten, F., Johansen, A. M., Naesseth, C. A., Kirkpatrick, B., Schön, T. B., Aston, J., and Bouchard-Côté, A. (2014). Divide-and-conquer with sequential Monte Carlo. *arXiv preprint arXiv:1406.4993*.
- Minka, T. (2005). Divergence measures and message passing. Technical report, Microsoft Research.
- Murray, L. M., Lee, A., and Jacob, P. E. (2014). Parallel resampling in the particle filter. *arXiv preprint arXiv:1301.4019*.
- Naesseth, C. A., Lindsten, F., and Schön, T. B. (2014). Sequential Monte Carlo for graphical models. In *Advances in Neural Information Processing Systems 27*.
- Naesseth, C. A., Lindsten, F., and Schön, T. B. (2015). Towards automated sequential Monte Carlo for probabilistic Graphical Models. In *NIPS Workshop on Black Box Learning and Inference*.
- Paige, B., Wood, F., Doucet, A., and Teh, Y. W. (2014). Asynchronous anytime sequential Monte Carlo. In *Advances in Neural Information Processing Systems 27*, pages 3410–3418.
- Pearl, J. and Russell, S. (1998). *Bayesian networks*. Computer Science Department, University of California.
- Pitt, M. K. and Shephard, N. (1999). Filtering via simulation: auxiliary particle filter. *Journal of the American Statistical Association*, 94:590–599.
- Ritchie, D., Mildenhall, B., Goodman, N. D., and Hanrahan, P. (2015). Controlling procedural modeling programs with stochastically-ordered sequential Monte Carlo. *ACM Transactions on Graphics (TOG)*, 34(4):105.
- Stuhlmüller, A., Taylor, J., and Goodman, N. (2013). Learning stochastic inverses. In *Advances in Neural Information Processing Systems 26*, pages 3048–3056.
- Todeschini, A., Caron, F., Fuentes, M., Legrand, P., and Del Moral, P. (2014). Biips: software for Bayesian inference with interacting particle systems. *arXiv preprint arXiv:1412.3779*.
- Uria, B., Murray, I., and Larochelle, H. (2013). RNADE: The real-valued neural autoregressive density-estimator. In *Advances in Neural Information Processing Systems*, pages 2175–2183.
- Van Der Merwe, R., Doucet, A., De Freitas, N., and Wan, E. (2000). The unscented particle filter. In *Advances in Neural Information Processing Systems*, pages 584–590.
- Wood, F., van de Meent, J. W., and Mansinghka, V. (2014). A new approach to probabilistic programming inference. In *Proceedings of the 17th International conference on Artificial Intelligence and Statistics*.