
Resampled Belief Networks for Variational Inference

Ramki Gummadi*
gsrk84@gmail.com

1 Introduction

We consider the problem of learning representations using a combination architecture consisting of generative/recognition networks, where the goal is to simultaneously achieve rich representations along with fast inference/training, similar in spirit to [1, 2, 3, 4, 5, 6]. The use of a directed belief network for the recognition part helps in training the generative model, but is also useful purely as a tool for efficient inference (e.g. [7, 8]). For a generative model density p parametrized by θ , and any other density q (e.g. the recognition density, parametrized by ϕ), conditioning on the evidence \mathbf{v} , we have the following equations (reviewed in Section 7.1):

$$\log p(\mathbf{v}) = \mathbb{E}_{q(\mathbf{h}|\mathbf{v})} \left[\log \frac{p(\mathbf{h}, \mathbf{v})}{q(\mathbf{h}|\mathbf{v})} \right] + \text{KL}(q(\mathbf{h}|\mathbf{v})||p(\mathbf{h}|\mathbf{v}))$$

$$\nabla_{\theta} \log p(\mathbf{v}) = \mathbb{E}_{p(\mathbf{h}|\mathbf{v})} [\nabla_{\theta} \log p(\mathbf{h}, \mathbf{v})] = \mathbb{E}_{q(\mathbf{h}|\mathbf{v})} [\nabla_{\theta} \log p(\mathbf{h}, \mathbf{v})] + \nabla_{\theta} \text{KL}(q(\mathbf{h}|\mathbf{v})||p(\mathbf{h}|\mathbf{v}))$$

Although the generative and the recognition network each individually contains a representation of the data, the task of jointly training them may be viewed as an approximation to training the generative model alone. Variational methods are able to drop the second term in the above gradient for θ by instead estimating $\nabla_{\phi} \text{KL}(q(\mathbf{h}|\mathbf{v})||p(\mathbf{h}|\mathbf{v}))$ and then performing gradient ascent on an expanded collection of parameters (θ, ϕ) . However, the path taken by θ during gradient ascent on the *evidence lower bound* in variational methods (which corresponds to dropping the KL divergence) clearly depends on the class over which q is optimized. In mean field methods [9] that use independence assumptions, the parameters ϕ need to be recomputed from scratch for each inference instance. By contrast, using a directed recognition network with parametrized edge weights allows learning an inverse factorization of the generative model directly. However: (1) it may not necessarily contain the density implied for inference; and (2) training the recognition model parameters is a challenging problem, even with a fixed θ .

We propose a class of parametrized densities called Resampled Belief Networks (RBN), which are: (1) at least theoretically, capable of representing the exact inverse and (2) introduce an additional lever to control the variance of gradient estimates. An RBN augments the recognition directed belief network with a resampling scheme designed to tighten the gap between the recognition and generative densities. In Section 2, we characterize the density achieved by an RBN in terms of the constituent directed network and the resampling rules. Using this characterization, in Sections 3 and 4 we develop a scheme to break down a given recognition/generative model into corresponding resampling constraints. Another motivation for this approach is to achieve tradeoff points along a continuum between monte carlo approaches with perfect accuracy but impractical computational expense, and variational approximations with much better computational expense.

2 Resampled Belief Networks (RBN)

Let R be a directed belief network on $\mathbf{x} = (x_1, \dots, x_n)$ that samples from x_1 to x_n using: $r(\mathbf{x}) = \prod_{i=1}^n r_i(x_i|\mathbf{x}_{pa_r(i)})$. By an RBN, we refer to the sampler defined in Algorithm 1,

*This research was primarily done while the author was with the Department of Computing Science, University of Alberta. Current homepage: <http://ramki-gummadi.github.io/>

by augmenting R with (i) *resampling constraints* specified using a collection of constraint functions, $C = \{C_1, \dots, C_n\}$, where each constraint¹, $C_i : \mathbf{x}_{[i]} \mapsto (0, 1]$; and (ii) *jump rules* $\Delta = \{\Delta_1, \dots, \Delta_n\}$ where $\Delta_i \in \{0, \dots, i-1\}$ specifies the node from which resampling needs to be done when constraint i fails.

Algorithm 1 Sampling definition of RBN(R, C, Δ).

- 1: $i = 1$.
 - 2: **while** $i \leq n$ **do**
 - 3: $x_i \leftarrow$ sample from $r_i(\cdot | \mathbf{x}_{pa_r(i)})$.
 - 4: $u \sim U[0, 1]$.
 - 5: **if** $u > C_i(\mathbf{x})$ **then**
 - 6: $i \leftarrow \Delta_i$.
 - 7: $i \leftarrow i + 1$.
 - 8: Output sample $\mathbf{x} = (x_1, \dots, x_n)$.
-

Theorem 1 (Density characterization). *Let $Z_i(\mathbf{x}) \triangleq \mathbb{E}[C_i(\mathbf{x}) | (x_1, \dots, x_{\Delta_i})]$, where the expectation is with respect to the density induced on (x_1, \dots, x_i) prior to constraint i . Then, the sampling density of Algorithm 1 satisfies $q(\mathbf{x}) = r(\mathbf{x}) \frac{\prod_i C_i(\mathbf{x})}{\prod_i Z_i(\mathbf{x})}$.*

Note that Z_i denotes the probability that a randomly generated sample is accepted by constraint i as a function of the sampled nodes prior to the resampling point.

Corollary 2. *If $Z_i(x_1, \dots, x_{\Delta_i}) = Z_i$, i.e. independent of its arguments, then $q(\mathbf{x}) = \frac{1}{Z} r(\mathbf{x}) \prod_i C_i(\mathbf{x})$, where Z is the appropriate normalizing constant. A useful special case is that of always resampling from the start, i.e. $\Delta_i = 0 \forall i$.*

3 Generative/Recognition Network architecture with RBN

Assume all variables are binary and let $p(\mathbf{x})$ denote the distribution implied by a generative model, parametrized by the vector θ . Let the recognition model contain a directed belief network with distribution $r(\mathbf{x})$, whose parameters are contained in the vector ϕ . The energy functions for relevant distributions will be composed from $a_i(\mathbf{x})$ and $b_i(\mathbf{x})$ defined below, where $pa(i)$ refers to the generative model.

$$a_i(\mathbf{x}) = \begin{cases} \sum_{j \in pa(i)} \theta_{ji} x_j & \text{if } x_i = 1 \\ -\sum_{j \in pa(i)} \theta_{ji} x_j & \text{if } x_i = 0 \end{cases}, \quad b_i(\mathbf{x}) = \begin{cases} \sum_{j < i} \phi_{ji} x_j & \text{if } x_i = 1 \\ -\sum_{j < i} \phi_{ji} x_j & \text{if } x_i = 0 \end{cases} \quad (1)$$

Using the shorthand operator $f^+(\mathbf{x}) \triangleq \log(1 + e^{f(\mathbf{x})})$, we will rewrite the usual expression for the distributions as follows.

$$r(\mathbf{x}) = \prod_{i \geq 0} \frac{e^{-(\sum_{j < i} \phi_{ji} x_j) x_i}}{1 + e^{-(\sum_{j < i} \phi_{ji} x_j)}} = \prod_{i \geq 0} e^{-b_i^+(\mathbf{x})}, \quad \text{and } p(\mathbf{x}) = \frac{1}{Z_p} \prod_i e^{-a_i^+(\mathbf{x})} \quad (2)$$

The identity is easily verified by considering $x_i \in \{0, 1\}$ separately.² Let the constraint functions be defined by their negative log likelihoods, $c_j(\mathbf{x}) \geq 0$ over the constraint index j , i.e., $C_j(\mathbf{x}) = e^{-c_j(\mathbf{x})}$. For now, we will only consider jump rules with $\Delta_i = 0 \forall i$, although this is not necessary to satisfy Equation (3) as long as Corollary 2 is valid. From Theorem 1, the recognition network samples from:

$$q(\mathbf{x}) = \frac{1}{Z_q} \prod_{i \geq 0} e^{-b_i^+(\mathbf{x})} \prod_j e^{-c_j(\mathbf{x})} \triangleq \frac{\gamma_q(\mathbf{x})}{Z_q} \quad (3)$$

¹Let $[i] \triangleq \{1, 2, \dots, i\}$. Furthermore, we abuse notation and write the element, rather than the set which is being mapped, for clarity. Note that C_i are functions that depend only on the first i variables of \mathbf{x} . However, for reducing notation clutter, we use $C_i(\mathbf{x})$ with the implicit recognition that the components of \mathbf{x} outside $[i]$ are irrelevant.

²Note that $Z_p = 1$, but we write it in the above equation explicitly to reuse the equation with separate hidden and visible variables. For that purpose, the index i may be allowed to have terms beyond the indices of the sampling variables.

Note that rejection sampling is analogous to using a single constraint function, $c(\mathbf{x}) = c + \sum_i a_i^+(\mathbf{x}) - \sum_{i \geq 0} b_i^+(\mathbf{x})$, where c is any constant chosen large enough to ensure that $c(\mathbf{x}) \geq 0 \forall \mathbf{x}$. An interpretation of the subclass of jump rules for which $\Delta_i = 0$, is to view them as approximately breaking down a large rejection event into smaller sequential pre-emptive steps, thereby improving the sampling efficiency, possibly exponentially.

Proposition 3.

$$\text{Let: } \mathcal{A}(\mathbf{x}) \triangleq \sum_i a_i^+(\mathbf{x}) - \sum_{i \geq 0} b_i^+(\mathbf{x}) - \sum_j c_j(\mathbf{x}), \text{ and } \bar{\mathcal{A}}(\mathbf{x}) \triangleq \mathcal{A}(\mathbf{x}) - \mathbb{E}_q[\mathcal{A}(\mathbf{x})] \quad (4)$$

$$\text{Then, } \text{KL}(q||p) = \log \left(\mathbb{E}_q \left[e^{-\bar{\mathcal{A}}(\mathbf{x})} \right] \right), \text{ and } \nabla_\phi \text{KL}(q||p) = \text{COV}(\mathcal{A}(\mathbf{x}), \nabla_\phi \mathcal{A}(\mathbf{x})) \quad (5)$$

Furthermore, when $\text{KL}(q||p) \approx 0$, an approximation to the KL divergence is $\text{KL}(q||p) \approx \text{VAR}(\mathcal{A}(\mathbf{x}))/2$.

The proof is provided via Theorem 7 in the Appendix. In fact, it is worth noting that Theorem 7 is stated for arbitrary unnormalized distributions and not specifically for RBNs. The claim about the approximation is justified by making use of the second order Taylor approximation $e^{-x} \approx 1 - x + x^2/2$ in Equation (5). While Equation (5) suggests the importance of controlling the variance of $\mathcal{A}(\mathbf{x})$ to estimate gradients for the recognition network, Proposition 3 also suggests that achieving a low variance for $\mathcal{A}(\mathbf{x})$ would require having access to a sampler with an approximately accurate inferential density, implying a chicken and egg situation. This justifies the claim made in the introduction that RBNs contain a lever for variance control of the gradient estimates.

4 Derivation of Resampling Constraints from the Generative Model

In this Section, we use the above framework to propose one possible general design for the resampling constraints without making any assumptions on the generative model. The basic observation underlying this proposal is that, under the target sampling density $p(\mathbf{x})$, the functions $a_i^+(\mathbf{x})$ have a uniform concentration bound over all possible parameters. This property, formalized in Lemma 4, therefore suggests that whenever a generated sample violates this concentration property by an appropriately large margin, it could be considered a credible signal to resample.

Lemma 4. $a_i^+(\mathbf{x}) < 0.7 + 0.75k$ with probability at least $1 - 1/k^2$, uniformly across θ, i and sigmoid generative sample densities $p(\mathbf{x})$ defined in Equation 2 (with $Z_p = 1$).

The proof is provided in the appendix. Let: $S_t = \{i : pa(i) \in [t] \text{ and } pa(i) \notin [t - 1]\}$, for $t \in [n]$. That is, S_t represents the constraints that are revealed at step t in the recognition network. It is helpful to rewrite Equation (4) as:

$$\mathcal{A}(\mathbf{x}) = \sum_{t=1}^n A_t(\mathbf{x}) - \sum_j c_j(\mathbf{x}), \text{ where } A_t(\mathbf{x}) \triangleq \sum_{i \in S_t} a_i^+(\mathbf{x}) - b_t^+(\mathbf{x}) \quad (6)$$

For an arbitrarily general factorization of $p(\mathbf{x})$ with no structure, S_t may be non-empty only when $t = n$, but it could be subject to a more uniform “revelation schedule” in structured factorizations, e.g. layered networks ([3]). The concentration property for individual $a_i^+(\mathbf{x})$ in Lemma 4 can be generalized to the functions $A_t(\mathbf{x})$, even though the different $a_i(\mathbf{x})$ are correlated, to obtain appropriate bounds that only depend on the size of S_t . This motivates the class of constraint functions in Equation (7), where $\{T_t : t \in [n]\}$ may be considered as additional parameters complementary to the edge weights ϕ_{lm} in determining the recognition density q .

$$c_t(\mathbf{x}) = (A_t(\mathbf{x}) - T_t)^+ \quad \forall t \in [n] \quad (7)$$

This implies: $\mathcal{A}(\mathbf{x}) = \sum_{t=1}^n (A_t(\mathbf{x}) - (A_t(\mathbf{x}) - T_t)^+)$. From this equation, it is clear that setting T_t small enough will ensure that $\text{Var}(\mathcal{A}(\mathbf{x})) \approx 0$, which in turn implies that $\text{KL}(q||p) \approx 0$. However, this could be infeasible due to an untenable sampling efficiency. By contrast, setting T_t large disables all constraints, and the sampling density of the RBN remains unaltered from R .

5 Experimental Observations

We run numerical experiments using two synthetic binary 100 node sigmoid belief network generative/recognition architectures, (R_i, G_i) for $i \in \{\text{rev}, \text{nr}\}$, where R_i is part of the recognition density and G_i is its generative counterpart. Our initial focus is not directly on training (R_i, G_i) , but to first study the effects of introducing the proposed resampling framework on top of the standard architecture. G_{rev} is chosen to correspond to a reversible markov chain, due to its property that its exact inverse factorization has identical parameters as G_i . However, the exclusively one step dependencies might be expected to result in weak constraints across nodes, suggesting that this may be a pathological case for demonstrating any potential benefits of resampling. R_{rev} was fixed to be a small perturbation from the perfect inverse factorization. By contrast, G_{nr} had randomly sampled edge weights from each node upto its 10 immediate parents. R_{nr} had full connectivity with randomly sampled edge weights. The choice of the thresholds T_t , left unspecified in Section 4, is chosen according to the following heuristic for both experiments: First, for each $t \in [n]$, we form Monte Carlo estimates for the mean, $\hat{\mu}_t$ and standard deviation, $\hat{\sigma}_t$, for samples of $A_t(\mathbf{x})$ defined in Equation (6) under G_i . Given a uniform (across t) control parameter, called `gap`, we then choose $T_t = \hat{\mu}_t + \text{gap} * \hat{\sigma}_t$. Finally, `gap` is adjusted according to a multiplicative rate control heuristic that responds to the sampling efficiency. We find that the control variable `gap` is more instrumental in the case of `nrev`. The following three statistics are shown in the appendix in Figure 1 across the trajectory of the `gap` variable for the resulting RBNs, on each of `rev` and `nrev`: (1) NLL (2) Sample efficiency (3) $\text{VAR}(\mathcal{A}(\mathbf{x}))/2$ (which may be considered a proxy for the KL distance achieved as explained in Proposition 3).

6 Amortization of Inference and Training Equations

In this section, we consider the problem of training the RBN. Split $\mathbf{x} = (\mathbf{v}, \mathbf{h})$, where $\mathbf{v} = (x_1, \dots, x_i)$ and $\mathbf{h} = (x_{i+1}, \dots, x_n)$ represent the visible/hidden parts, and restrict ϕ to those parameters that only affect the conditional density of \mathbf{h} given \mathbf{v} . The gradient in Equation (5) can be rewritten using the law of total covariance as:

$$\nabla_{\phi} \text{KL}(q||p) = \mathbb{E}_{\mathbf{v}} [\text{COV}_{\mathbf{h}}(\mathcal{A}(\mathbf{x}), \nabla_{\phi} \log \gamma_q(\mathbf{x}))] + \text{COV}_{\mathbf{v}}(\mathbb{E}_{\mathbf{h}}[\mathcal{A}(\mathbf{x})], \mathbb{E}_{\mathbf{h}}[\nabla_{\phi} \log \gamma_q(\mathbf{x})]) \quad (8)$$

In the above equation, subscripts of operators refer to the variables over which we are not conditioning, i.e. considering as random variables. Using arguments similar to the proof of Theorem 7, we can write $\mathbb{E}_{\mathbf{h}}[\nabla_{\phi} \log \gamma_q(\mathbf{v}, \mathbf{h})] = \nabla_{\phi} \log Z_{q_{\mathbf{v}}}$, where $Z_{q_{\mathbf{v}}} = \sum_{\mathbf{h}} \gamma_q(\mathbf{v}, \mathbf{h})$. By defining $q_{\mathbf{v}}, p_{\mathbf{v}}$ as the conditional densities of q and p on \mathbf{h} given \mathbf{v} , and applying $\nabla_{\phi} \text{KL}(q_{\mathbf{v}}||p_{\mathbf{v}}) = \text{COV}_{\mathbf{h}}(\mathcal{A}(\mathbf{v}, \mathbf{h}), \nabla_{\phi} \log \gamma_q(\mathbf{v}, \mathbf{h}))$, we may write Equation (8) as:

$$\nabla_{\phi} \text{KL}(q||p) = \mathbb{E}_{\mathbf{v}} [\nabla_{\phi} \text{KL}(q_{\mathbf{v}}||p_{\mathbf{v}})] + \text{COV}_{\mathbf{v}}(\mathbb{E}_{\mathbf{h}}[\mathcal{A}(\mathbf{x})], \nabla_{\phi} \log Z_{q_{\mathbf{v}}}) \quad (9)$$

From the above equation, we see that $Z_{q_{\mathbf{v}}}$ being independent of \mathbf{v} (can still depend on i or the edge weights) is a sufficient condition to equate the second term is zero. In this case, the parameter training can be averaged out over different inference instances. For an RBN, $Z_{q_{\mathbf{v}}}$ is equal to the probability of acceptance of a sample \mathbf{x} starting from \mathbf{v} , provided that $\Delta_j = i \forall j > i$. Therefore, any node in the network where the probability of acceptance for samples generated from that point is independent of the values of the nodes that were sampled previously can serve as an amortization point for training the parameters of the recognition network. Since the acceptance probabilities in an RBN can be controlled via threshold parameters, this suggests one possible way to exploit this observation in practice. Theorem 5 provides a summary of the derivations for the gradients of the training parameters in an RBN.

Theorem 5. Let ϕ_{lm} be the weight parameter connecting node x_l to node x_m ($l < m$) in the recognition network, R . Let T_m denote the threshold parameter in $c_m(\mathbf{x})$ according to Equation (7) for the RBN. With $\sigma(\mathbf{x}) \triangleq (1 + e^{-\mathbf{x}})^{-1}$, we get the following expressions for gradients, where we define $\mu_m(\mathbf{x})$ as the conditional expected value of x_m for the belief network R , given all its parents.

$$-\frac{\partial \text{KL}(q||p)}{\partial \phi_{lm}} = \text{COV}(\mathcal{A}(\mathbf{x}), (1 - \mu_m(\mathbf{x}))x_l x_m \sigma(T_m - A_m(\mathbf{x})))$$

$$\frac{\partial \text{KL}(q||p)}{\partial T_m} = \text{COV}(\mathcal{A}(\mathbf{x}), \sigma(A_m(\mathbf{x}) - T_m))$$

An interesting direction for future work is to evaluate simultaneous training of edge weight parameters alongside learning the resampling thresholds and comparing it with the heuristic learning rule that was used in Section 5.

7 Appendix

This Section provides proof details, as well as a figure reporting some numerical observations in experiments.

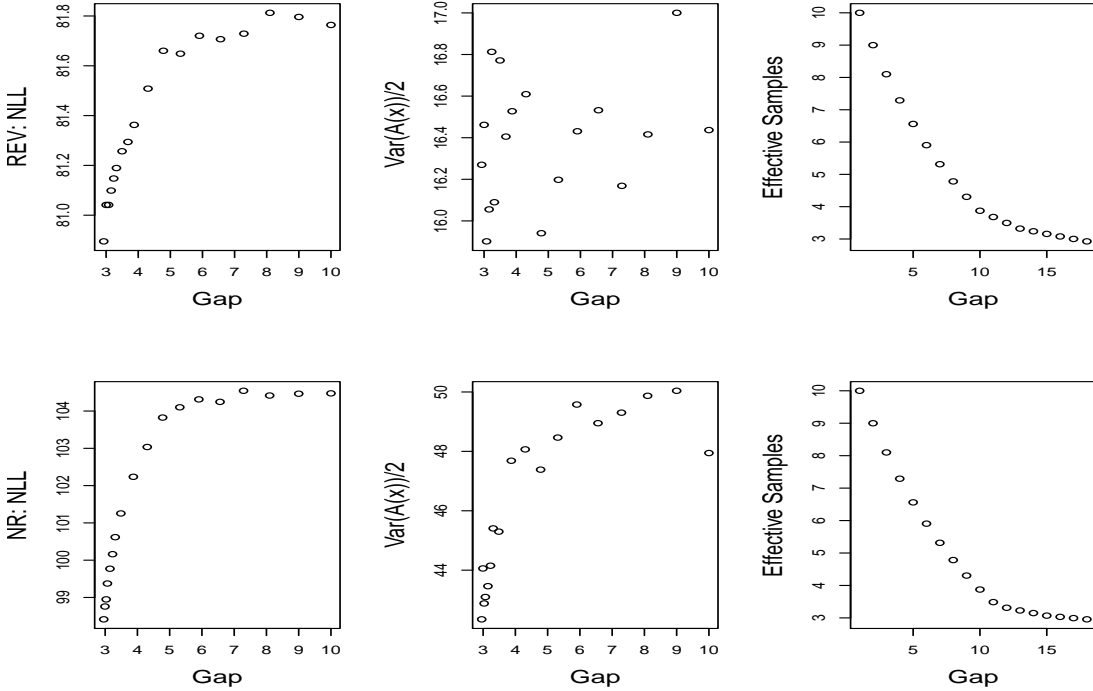


Figure 1: The top row figures show the statistics for a reversible Markov Chain on 100 binary nodes, for which the conditionals $p(x_{i+1}|\mathbf{x}_{1:i}) = p(x_{i+1}|x_i)$, correspond to the transition matrix $\begin{bmatrix} 0.7 & 0.3 \\ 0.2 & 0.8 \end{bmatrix}$. The prior for the first node was chosen to achieve reversibility, i.e. the stationary distribution. The bottom row depicts results from an experiment where the generative network was fixed using a distribution where each node has a conditional that depends on the previous ten sampled nodes, where the corresponding edge weights were chosen by sampling from a standard normal. The gap variable refers to a uniform control parameter for the threshold parameters across all nodes according to the mechanism described in Section 5. Each depicted point is an average over 10,000 samples of the corresponding RBN resulting from the threshold rule defined by the gap variable. NLL refers to the negative log likelihood of samples generated by the RBN evaluated with respect to the corresponding generative density. The effective samples depict the factor overhead that results from resampling. All resampling rules implemented correspond to $\Delta_i = 0 \forall i$.

Theorem (Density characterization). *Let $Z_i(\mathbf{x}) \triangleq \mathbb{E}[C_i(\mathbf{x})|(x_1, \dots, x_{\Delta_i})]$, where the expectation is with respect to the density induced on (x_1, \dots, x_i) prior to constraint i . Note that Z_i denotes the probability that constraint i fails, as a function of the nodes prior to the resampling point. Then, the sampling density of Algorithm 1 satisfies $q(\mathbf{x}) = r(\mathbf{x}) \frac{\prod_i C_i(\mathbf{x})}{\prod_i Z_i(\mathbf{x})}$.*

Proof. [Proof of Theorem 1] Wherever convenient, we will use $\mathbf{x}_{1:i}$ to denote (x_1, \dots, x_i) . Let $q_k(x_{1:k})$ be the sampling density of $\mathbf{x}_{1:k}$, at the point where $i = k$, just prior to Step 7. Therefore,

the density we are evaluating, $q(\mathbf{x})$, is equal to $q_n(\mathbf{x})$. Note that $q_k(\mathbf{x}_{1:k})$ is not necessarily equal to the marginal density of q_n on $\mathbf{x}_{1:k}$. Let $q_k(\mathbf{x}_{1:i})$ denote the marginal density of q_k on $\mathbf{x}_{1:i}$, and let $q_k(\mathbf{x}_{i+1:k}|\mathbf{x}_{1:i})$ denote the appropriate conditional density. The proof will use induction on n . The base case is trivial, and assume that the induction hypothesis holds for $n - 1$. Let $\text{accept} \triangleq \mathbb{1}\{u \leq C_n(\mathbf{x})\}$ in Step 5, so that $P(\text{accept}|X) = \mathbb{E}[C_n(\mathbf{x})|X]$ for any random variable X . Conditioning on $\mathbf{x}_{1:\Delta_n}$:

$$\begin{aligned} q_n((x_{\Delta_n+1}, \dots, x_n)|\mathbf{x}_{1:\Delta_n}) &= P((x_{\Delta_n+1}, \dots, x_n)|\text{accept}, \mathbf{x}_{1:\Delta_n}) \\ &= P(((x_{\Delta_n+1}, \dots, x_n), \text{accept})|\mathbf{x}_{1:\Delta_n}) / P(\text{accept}|\mathbf{x}_{1:\Delta_n}) \\ &= q_{n-1}((x_{\Delta_n+1}, \dots, x_{n-1})|\mathbf{x}_{1:\Delta_n})r(x_n|\mathbf{x}_{1:n-1})C_n(\mathbf{x})/Z_n(\mathbf{x}_{1:\Delta_n}) \end{aligned}$$

Noting that the marginals of q_n and q_{n-1} are identical on $\mathbf{x}_{1:\Delta_n}$,

$$\begin{aligned} q_n(x_{1:n}) &= q_n(\mathbf{x}_{1:\Delta_n})q_n((x_{\Delta_n+1}, \dots, x_n)|\mathbf{x}_{1:\Delta_n}) \\ &= q_{n-1}(\mathbf{x}_{1:\Delta_n})q_{n-1}((x_{\Delta_n+1}, \dots, x_{n-1})|\mathbf{x}_{1:\Delta_n})r(x_n|\mathbf{x}_{1:n-1})C_n(\mathbf{x})/Z_n(\mathbf{x}_{1:\Delta_n}) \\ &= q_{n-1}(\mathbf{x}_{1:n-1})r(x_n|\mathbf{x}_{1:n-1})C_n(\mathbf{x})/Z_n(\mathbf{x}_{1:\Delta_n}) \end{aligned}$$

The proof follows by appealing to the induction hypothesis in the last equation. \square

Remark 6. The constraints were defined over the index set $[n]$ for clarity of presentation. However, the characterization identified also generalizes to an arbitrary index set, with $C_j : \mathbf{x}_{A_j} \mapsto (0, 1]$. In this case, Step 5 is a loop over $\{j : A_j \subseteq [i] \text{ and } A_j \not\subseteq [i-1]\}$. A natural resampling rule to consider here would be: $\Delta_j = \min\{A_j\} - 1$, although the characterization remains valid regardless.

Theorem 7. Suppose $p(\mathbf{x}) = \gamma_p(\mathbf{x})/Z_p$ and $q(\mathbf{x}) = \gamma_q(\mathbf{x})/Z_q$, where only q depends on ϕ . Let $\mathcal{A}(\mathbf{x}) \triangleq \log \gamma_q(\mathbf{x}) - \log \gamma_p(\mathbf{x})$; and $\bar{\mathcal{A}}(\mathbf{x}) \triangleq \mathcal{A}(\mathbf{x}) - \mathbb{E}_q[\mathcal{A}(\mathbf{x})]$. Then:

$$\text{KL}(q||p) = \log \left(\mathbb{E}_q \left[e^{-\bar{\mathcal{A}}(\mathbf{x})} \right] \right) \quad (10)$$

$$\nabla_\phi \text{KL}(q||p) = \text{COV}_q(\mathcal{A}(\mathbf{x}), \nabla_\phi \mathcal{A}(\mathbf{x})) \quad (11)$$

where the subscript q refers to the density used to generate the samples of \mathbf{x} .

Proof.

$$\begin{aligned} \text{KL}(q||p) &= \log Z_p - \log Z_q + \mathbb{E}_q[\log \gamma_q(\mathbf{x}) - \log \gamma_p(\mathbf{x})] \\ &= \log \left(\sum_{\mathbf{x}} \frac{\gamma_p(\mathbf{x})}{Z_q} \right) + \mathbb{E}_q[\log \gamma_q(\mathbf{x}) - \log \gamma_p(\mathbf{x})] \\ &= \log \left(\mathbb{E}_q \left[\frac{\gamma_p(\mathbf{x})}{\gamma_q(\mathbf{x})} \right] \right) + \mathbb{E}_q[\log \gamma_q(\mathbf{x}) - \log \gamma_p(\mathbf{x})] \\ &= \log \left(\mathbb{E}_q \left[e^{-\mathcal{A}(\mathbf{x})} \right] \right) + \mathbb{E}_q[\mathcal{A}(\mathbf{x})] = \log \left(\mathbb{E}_q \left[e^{-\bar{\mathcal{A}}(\mathbf{x})} \right] \right) \end{aligned}$$

For the gradients, differentiating the first equality for $\text{KL}(q||p)$, we get: $\nabla_\phi \text{KL}(q||p) = D_1 - D_2 - D_3$, where $D_1 = \nabla_\phi \mathbb{E}_q[\log \gamma_q(\mathbf{x})]$, $D_2 = \nabla_\phi \mathbb{E}_q[\log \gamma_p(\mathbf{x})]$, $D_3 = \nabla_\phi \log Z_q$.

$$\begin{aligned} D_1 &= \nabla_\phi \mathbb{E}_q[\log \gamma_q(\mathbf{x})] = \sum_{\mathbf{x}} \nabla_\phi [q(\mathbf{x}) \log \gamma_q(\mathbf{x})] \\ &= \sum_{\mathbf{x}} \left(\frac{q(\mathbf{x})}{\gamma_q(\mathbf{x})} \nabla_\phi \gamma_q(\mathbf{x}) + \log \gamma_q(\mathbf{x}) \nabla_\phi q(\mathbf{x}) \right) \\ &= \frac{1}{Z_q} \nabla_\phi Z_q + \sum_{\mathbf{x}} q(\mathbf{x}) \log \gamma_q(\mathbf{x}) \nabla_\phi \log q(\mathbf{x}) \\ &= D_3 + \mathbb{E}_q[\log \gamma_q(\mathbf{x}) \nabla_\phi \log q(\mathbf{x})] \end{aligned}$$

Similarly, $D_2 = \mathbb{E}_q[\log \gamma_p(\mathbf{x}) \nabla_\phi \log q(\mathbf{x})]$, which implies:

$$\nabla_\phi \text{KL}(q||p) = \mathbb{E}_q[(\log \gamma_q(\mathbf{x}) - \log \gamma_p(\mathbf{x})) \nabla_\phi \log q(\mathbf{x})]$$

Note that $\nabla_{\phi} \log q(\mathbf{x}) = \nabla_{\phi} \log \gamma_q(\mathbf{x}) - D_3 = \nabla_{\phi} \log \gamma_q(\mathbf{x}) - \mathbb{E}_q [\nabla_{\phi} \log \gamma_q(\mathbf{x})]$, where:

$$D_3 = \nabla_{\phi} \log Z_q = \frac{1}{Z_q} \sum_{\mathbf{x}} \nabla_{\phi} \gamma_q(\mathbf{x}) = \frac{1}{Z_q} \sum_{\mathbf{x}} \gamma_q(\mathbf{x}) \nabla_{\phi} \log \gamma_q(\mathbf{x}) = \mathbb{E}_q [\nabla_{\phi} \log \gamma_q(\mathbf{x})]$$

This implies:

$$\nabla_{\phi} \text{KL}(q||p) = \mathbb{E}_q \left[(\log \gamma_q(\mathbf{x}) - \log \gamma_p(\mathbf{x})) \overline{\nabla_{\phi} \log \gamma_q(\mathbf{x})} \right] \quad (12)$$

$$= \text{COV}(\log \gamma_q(\mathbf{x}) - \log \gamma_p(\mathbf{x}), \nabla_{\phi} \log \gamma_q(\mathbf{x})) \quad (13)$$

which completes the proof since $\nabla_{\phi} \log \gamma_q(\mathbf{x}) = \nabla_{\phi} \mathcal{A}(\mathbf{x})$. \square

Lemma. $a_i^+(\mathbf{x}) < 0.7 + 0.75k$ with probability at least $1 - 1/k^2$, uniformly across θ, i and sigmoid generative sample densities $p(\mathbf{x})$ defined in Equation 2 (with $Z_p = 1$).

Proof. [Proof of Lemma 4] Denote $\alpha_i \triangleq \left(\sum_{j \in \text{pa}(i)} \theta_{ji} x_j \right)$. Then, $\text{VAR}(a_i^+(\mathbf{x})) \leq \mathbb{E}[(a_i^+(\mathbf{x}))^2] = \mathbb{E}[\mathbb{E}[(a_i^+(\mathbf{x}))^2 | \alpha_i]] < 0.6$ from Lemma 8. Similarly, $\mathbb{E}[a_i^+(\mathbf{x}) | \alpha_i] \leq \ln 2$, again from Lemma 8. The claim then follows as an application of Chebychev's inequality. It is also feasible to obtain tighter concentration bounds using exponential versions of the inequality. \square

Lemma 8. Suppose $X = [\alpha(\mathbb{1}\{x=1\} - \mathbb{1}\{x=0\})]^+$, for any constant α , where x has the sigmoid distribution with probability $\frac{e^{-\alpha x}}{1+e^{-\alpha}}$, for $x \in \{0, 1\}$ and $a^+ \triangleq \log(1 + e^{\alpha})$. Then: $0 \leq \mathbb{E}[X] < \log 2 \approx 0.693$ and $\mathbb{E}[X^2] < 0.6$ and $\text{VAR}(X) < 0.5$ for any α .

Proof. The proof of this claim can be obtained by explicitly computing the below functions of α for $k = 1, 2$:

$$\mathbb{E}[X^k] = \frac{\log^k(1 + e^{\alpha})}{1 + e^{\alpha}} + \frac{\log^k(1 + e^{-\alpha})}{1 + e^{-\alpha}}$$

Note also, that arbitrarily high moments of X have bounds that are independent of α . \square

Theorem. Let ϕ_{lm} be the weight parameter connecting node x_l to node x_m ($l < m$) in the recognition network, R . Let T_m denote the threshold parameter in $c_m(\mathbf{x})$ according to Equation 7 for the RBN. With $\sigma(\mathbf{x}) \triangleq (1 + e^{-\mathbf{x}})^{-1}$, we get the following expressions for gradients, where we define $\mu_m(\mathbf{x})$ as the conditional expected value of x_m for the belief network R , given all its parents.

$$\begin{aligned} -\frac{\partial \text{KL}(q||p)}{\partial \phi_{lm}} &= \text{COV}(\mathcal{A}(\mathbf{x}), (1 - \mu_m(\mathbf{x}))x_l x_m \sigma(T_m - A_m(\mathbf{x}))) \\ \frac{\partial \text{KL}(q||p)}{\partial T_m} &= \text{COV}(\mathcal{A}(\mathbf{x}), \sigma(A_m(\mathbf{x}) - T_m)) \end{aligned}$$

Proof. [Proof of Theorem 5]

$$\begin{aligned} -\nabla_{\phi} \log \gamma_q(\mathbf{x}) &= \nabla_{\phi} \left(\sum_{i \geq 0} b_i^+(\mathbf{x}) + \sum_j c_j(\mathbf{x}) \right) \\ &= \sum_{i \geq 0} \sigma(b_i(\mathbf{x})) \nabla_{\phi} b_i(\mathbf{x}) + \sum_j \nabla_{\phi} c_j(\mathbf{x}), \text{ where } \sigma(\mathbf{x}) = (1 + e^{-\mathbf{x}})^{-1} \end{aligned}$$

Consider the parameter ϕ_{lm} , which is the weight parameter connecting node x_l to node x_m in the recognition network, where $l < m$. The only non-zero term in the first summation is for $i = m$. Recalling Equation (1), we obtain:

$$-\frac{\partial \log \gamma_q(\mathbf{x})}{\partial \phi_{lm}} = \sigma(b_m(\mathbf{x}))x_l x_m + \sum_j \frac{\partial c_j(\mathbf{x})}{\partial \phi_{lm}} = (1 - \mu_m(\mathbf{x}))x_l x_m + \sum_j \frac{\partial c_j(\mathbf{x})}{\partial \phi_{lm}} \quad (14)$$

where, we define $\mu_m(\mathbf{x})$ as the conditional expected value of x_m given all its parents in the recognition network, R , so that $\sigma(b_m(\mathbf{x})) = 1 - \mu_m(\mathbf{x})$.

Referring back to Equation (14), we see from Equation (7) that $\frac{\partial c_j(\mathbf{x})}{\partial \phi_{lm}} \neq 0$ only for $j = m$, which is:

$$\frac{\partial c_m(\mathbf{x})}{\partial \phi_{lm}} = -\sigma(A_m(\mathbf{x}) - T_m) \frac{\partial b_m^+(\mathbf{x})}{\partial \phi_{lm}} = -\sigma(A_m(\mathbf{x}) - T_m) (1 - \mu_m(\mathbf{x})) x_l x_m$$

Using this to simplify Equation (14), we get:

$$-\frac{\partial \log \gamma_q(\mathbf{x})}{\partial \phi_{lm}} = (1 - \mu_m(\mathbf{x})) x_l x_m (1 - \sigma(A_m(\mathbf{x}) - T_m)) = (1 - \mu_m(\mathbf{x})) x_l x_m \sigma(T_m - A_m(\mathbf{x})) \quad (15)$$

The corresponding equation for T_m is:

$$\frac{\partial \log \gamma_q(\mathbf{x})}{\partial T_m} = \sigma(A_m(\mathbf{x}) - T_m) \quad (16)$$

□

7.1 The Variational Equation and its Gradients

This subsection contains a review of the variational equation and its gradients. Let \mathbf{v}, \mathbf{h} be visible/hidden parts of a random vector. Let θ be a scalar parameter for their joint distribution, $p(\mathbf{v}, \mathbf{h})$.

$$\begin{aligned} \frac{\partial \log p(\mathbf{v})}{\partial \theta} &= \frac{1}{p(\mathbf{v})} \frac{\partial p(\mathbf{v})}{\partial \theta} = \frac{1}{p(\mathbf{v})} \sum_{\mathbf{h}} \frac{\partial p(\mathbf{v}, \mathbf{h})}{\partial \theta} = \frac{1}{p(\mathbf{v})} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial \log p(\mathbf{v}, \mathbf{h})}{\partial \theta} \\ &= \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial \log p(\mathbf{v}, \mathbf{h})}{\partial \theta} = \mathbb{E}_{p(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log p(\mathbf{v}, \mathbf{h})}{\partial \theta} \right] \end{aligned}$$

Suppose, we use a tractable $q(\mathbf{h}|\mathbf{v})$ to compute the above expectation instead:

$$\frac{\partial \log p(\mathbf{v})}{\partial \theta} = \mathbb{E}_{p(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log p(\mathbf{v}, \mathbf{h})}{\partial \theta} \right] = \mathbb{E}_{q(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log p(\mathbf{v}, \mathbf{h})}{\partial \theta} \right] + \frac{\partial \text{KL}(q(\mathbf{h}|\mathbf{v})||p(\mathbf{h}|\mathbf{v}))}{\partial \theta} \quad (17)$$

Even if we solve the inference problem approximately accurately (i.e. with q close to p), we may be off the mark when attempting to improve the model because of the second term. To obtain the second equality of Equation (17), consider:

$$\begin{aligned} \log p(\mathbf{v}) &= \log p(\mathbf{v}, \mathbf{h}) - \log p(\mathbf{h}|\mathbf{v}) && \text{for any } \mathbf{h} \\ &= \sum_{\mathbf{h}} q(\mathbf{h}|\mathbf{v}) \log p(\mathbf{h}, \mathbf{v}) - \sum_{\mathbf{h}} q(\mathbf{h}|\mathbf{v}) \log p(\mathbf{h}|\mathbf{v}) && \text{for any distribution } q \text{ on } \mathbf{h} \\ &= \sum_{\mathbf{h}} q(\mathbf{h}|\mathbf{v}) \log \frac{p(\mathbf{h}, \mathbf{v})}{q(\mathbf{h}|\mathbf{v})} + \sum_{\mathbf{h}} q(\mathbf{h}|\mathbf{v}) \log \frac{q(\mathbf{h}|\mathbf{v})}{p(\mathbf{h}|\mathbf{v})} \\ &= \mathbb{E}_{q(\mathbf{h}|\mathbf{v})} \left[\log \frac{p(\mathbf{h}, \mathbf{v})}{q(\mathbf{h}|\mathbf{v})} \right] + \text{KL}(q(\mathbf{h}|\mathbf{v})||p(\mathbf{h}|\mathbf{v})) \end{aligned}$$

Denote the first term above as a function of p, q, \mathbf{v} where $\mathbb{H}(\cdot)$ denotes the Shannon entropy:

$$\mathcal{L}(p, q, \mathbf{v}) \triangleq \mathbb{E}_{q(\mathbf{h}|\mathbf{v})} \left[\log \frac{p(\mathbf{h}, \mathbf{v})}{q(\mathbf{h}|\mathbf{v})} \right] = \mathbb{E}_{q(\mathbf{h}|\mathbf{v})} [\log p(\mathbf{h}, \mathbf{v})] + \mathbb{H}(q(\mathbf{h}|\mathbf{v})) \quad (18)$$

Let ϕ be a parameter in the inferential/recognition density q . The variational approach to ML estimation involves gradient ascent on \mathcal{L} for the entire collection of θ and ϕ simultaneously:

$$\frac{\partial \mathcal{L}(p, q, \mathbf{v})}{\partial \theta} = \mathbb{E}_{q(\mathbf{h}|\mathbf{v})} \left[\frac{\partial \log p(\mathbf{h}, \mathbf{v})}{\partial \theta} \right] \quad \text{and} \quad \frac{\partial \mathcal{L}(p, q, \mathbf{v})}{\partial \phi} = -\frac{\partial \text{KL}_{\mathbf{v}}(q||p)}{\partial \phi} \quad (19)$$

To justify the gradient wrt ϕ , note that $\log p(\mathbf{v}) = \mathcal{L}(p, q, \mathbf{v}) + \text{KL}_{\mathbf{v}}(q||p)$ is independent of q . Therefore, we can interpret Equations (19) as compensation for dropping the term $\frac{\partial \text{KL}_{\mathbf{v}}(q||p)}{\partial \theta}$ in Equation (17) for gradient ascent on $\log p(\mathbf{v})$ wrt θ , by instead doing gradient ascent for an expanded variational collection that includes both θ and ϕ on $\mathcal{L}(p, q, \mathbf{v})$, by estimating the sensitivity of the KL distance with respect to the sampling density q instead.

8 Acknowledgements

I would like to thank anonymous reviewers for feedback and Prof. Dale Schuurmans for his support.

References

- [1] Peter Dayan, Geoffrey E Hinton, Radford M Neal, and Richard S Zemel. The helmholtz machine. *Neural computation*, 7(5):889–904, 1995.
- [2] Karol Gregor, Ivo Danihelka, Andriy Mnih, Charles Blundell, and Daan Wierstra. Deep autoregressive networks. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 1242–1250, 2014.
- [3] Karol Gregor, Andriy Mnih, and Daan Wierstra. Deep autoregressive networks. *arXiv preprint arXiv:1310.8499*, 2013.
- [4] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [5] Andriy Mnih and Karol Gregor. Neural variational inference and learning in belief networks. *arXiv preprint arXiv:1402.0030*, 2014.
- [6] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic back-propagation and variational inference in deep latent gaussian models. *arXiv preprint arXiv:1401.4082*, 2014.
- [7] Ross D Shachter and Mark Alan Peot. Simulation approaches to general probabilistic inference on belief networks. *arXiv preprint arXiv:1304.1526*, 2013.
- [8] Andreas Stuhlmüller, Jacob Taylor, and Noah Goodman. Learning stochastic inverses. In *Advances in neural information processing systems*, pages 3048–3056, 2013.
- [9] Martin J. Wainwright and Michael I. Jordan. Graphical models, exponential families, and variational inference. *Found. Trends Mach. Learn.*, 1(1-2):1–305, January 2008.