

### 18.7.1 Annealed Importance Sampling

In situations where  $D_{\text{KL}}(p_0||p_1)$  is large (i.e., where there is little overlap between  $p_0$  and  $p_1$ ), a strategy called **annealed importance sampling** (AIS) attempts to bridge the gap by introducing intermediate distributions (Jarzynski, 1997; Neal, 2001). Consider a sequence of distributions  $p_{\eta_0}, \dots, p_{\eta_n}$ , with  $0 = \eta_0 < \eta_1 < \dots < \eta_{n-1} < \eta_n = 1$  so that the first and last distributions in the sequence are  $p_0$  and  $p_1$ , respectively.

This approach enables us to estimate the partition function of a multimodal distribution defined over a high-dimensional space (such as the distribution defined by a trained RBM). We begin with a simpler model with a known partition function (such as an RBM with zeros for weights) and estimate the ratio between the two model's partition functions. The estimate of this ratio is based on the estimate of the ratios of a sequence of many similar distributions, such as the sequence of RBMs with weights interpolating between zero and the learned weights.

We can now write the ratio  $\frac{Z_1}{Z_0}$  as

$$\frac{Z_1}{Z_0} = \frac{Z_1}{Z_0} \frac{Z_{\eta_1}}{Z_{\eta_1}} \dots \frac{Z_{\eta_{n-1}}}{Z_{\eta_{n-1}}} \quad (18.47)$$

$$= \frac{Z_{\eta_1}}{Z_0} \frac{Z_{\eta_2}}{Z_{\eta_1}} \dots \frac{Z_{\eta_{n-1}}}{Z_{\eta_{n-2}}} \frac{Z_1}{Z_{\eta_{n-1}}} \quad (18.48)$$

$$= \prod_{j=0}^{n-1} \frac{Z_{\eta_{j+1}}}{Z_{\eta_j}}. \quad (18.49)$$

Provided the distributions  $p_{\eta_j}$  and  $p_{\eta_{j+1}}$ , for all  $0 \leq j \leq n-1$ , are sufficiently close, we can reliably estimate each of the factors  $\frac{Z_{\eta_{j+1}}}{Z_{\eta_j}}$  using simple importance sampling and then use these to obtain an estimate of  $\frac{Z_1}{Z_0}$ .

Where do these intermediate distributions come from? Just as the original proposal distribution  $p_0$  is a design choice, so is the sequence of distributions  $p_{\eta_1} \dots p_{\eta_{n-1}}$ . That is, it can be specifically constructed to suit the problem domain. One general purpose and popular choice for the intermediate distributions is to use the weighted geometric average of the target distribution  $p_1$  and the starting proposal distribution (for which the partition function is known)  $p_0$ :

$$p_{\eta_j} \propto p_1^{\eta_j} p_0^{1-\eta_j}. \quad (18.50)$$

In order to sample from these intermediate distributions, we define a series of Markov chain transition functions  $T_{\eta_j}(\mathbf{x}' | \mathbf{x})$  that define the conditional probability distribution of transitioning to  $\mathbf{x}'$  given we are currently at  $\mathbf{x}$ . The transition operator  $T_{\eta_j}(\mathbf{x}' | \mathbf{x})$  is defined to leave  $p_{\eta_j}(\mathbf{x})$  invariant:

$$p_{\eta_j}(\mathbf{x}) = \int p_{\eta_j}(\mathbf{x}') T_{\eta_j}(\mathbf{x} | \mathbf{x}') d\mathbf{x}'. \quad (18.51)$$

These transitions may be constructed as any Markov chain Monte Carlo method (e.g., Metropolis-Hastings, Gibbs), including methods involving multiple passes through all the random variables or other kinds of iterations.

The AIS sampling strategy is then to generate samples from  $p_0$  and use the transition operators to sequentially generate samples from the intermediate distributions until we arrive at samples from the target distribution  $p_1$ :

- for  $k = 1 \dots K$ 
  - Sample  $\mathbf{x}_{\eta_1}^{(k)} \sim p_0(\mathbf{x})$
  - Sample  $\mathbf{x}_{\eta_2}^{(k)} \sim T_{\eta_1}(\mathbf{x}_{\eta_2}^{(k)} \mid \mathbf{x}_{\eta_1}^{(k)})$
  - ...
  - Sample  $\mathbf{x}_{\eta_{m-1}}^{(k)} \sim T_{\eta_{m-2}}(\mathbf{x}_{\eta_{m-1}}^{(k)} \mid \mathbf{x}_{\eta_{m-2}}^{(k)})$
  - Sample  $\mathbf{x}_{\eta_m}^{(k)} \sim T_{\eta_{m-1}}(\mathbf{x}_{\eta_m}^{(k)} \mid \mathbf{x}_{\eta_{m-1}}^{(k)})$
- end

For sample  $k$ , we can derive the importance weight by chaining together the importance weights for the jumps between the intermediate distributions given in equation 18.49:

$$w^{(k)} = \frac{\tilde{p}_{\eta_1}(\mathbf{x}_{\eta_1}^{(k)}) \tilde{p}_{\eta_2}(\mathbf{x}_{\eta_2}^{(k)})}{\tilde{p}_0(\mathbf{x}_{\eta_1}^{(k)}) \tilde{p}_{\eta_1}(\mathbf{x}_{\eta_2}^{(k)})} \cdots \frac{\tilde{p}_1(\mathbf{x}_1^{(k)})}{\tilde{p}_{\eta_{m-1}}(\mathbf{x}_{\eta_m}^{(k)})}. \quad (18.52)$$

To avoid numerical issues such as overflow, it is probably best to compute  $\log w^{(k)}$  by adding and subtracting log probabilities, rather than computing  $w^{(k)}$  by multiplying and dividing probabilities.

With the sampling procedure thus defined and the importance weights given in equation 18.52, the estimate of the ratio of partition functions is given by:

$$\frac{Z_1}{Z_0} \approx \frac{1}{K} \sum_{k=1}^K w^{(k)}. \quad (18.53)$$

To verify that this procedure defines a valid importance sampling scheme, we can show (Neal, 2001) that the AIS procedure corresponds to simple importance sampling on an extended state space, with points sampled over the product space  $[\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{m-1}}, \mathbf{x}_1]$ . To do this, we define the distribution over the extended space as

$$\tilde{p}(\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{m-1}}, \mathbf{x}_1) \quad (18.54)$$

$$= \tilde{p}_1(\mathbf{x}_1) \tilde{T}_{\eta_{m-1}}(\mathbf{x}_{\eta_{m-1}} \mid \mathbf{x}_1) \tilde{T}_{\eta_{m-2}}(\mathbf{x}_{\eta_{m-2}} \mid \mathbf{x}_{\eta_{m-1}}) \cdots \tilde{T}_{\eta_1}(\mathbf{x}_{\eta_1} \mid \mathbf{x}_{\eta_2}), \quad (18.55)$$

where  $\tilde{T}_a$  is the reverse of the transition operator defined by  $T_a$  (via an application of Bayes' rule):

$$\tilde{T}_a(\mathbf{x}' | \mathbf{x}) = \frac{p_a(\mathbf{x}')}{p_a(\mathbf{x})} T_a(\mathbf{x} | \mathbf{x}') = \frac{\tilde{p}_a(\mathbf{x}')}{\tilde{p}_a(\mathbf{x})} T_a(\mathbf{x} | \mathbf{x}'). \quad (18.56)$$

Plugging the above into the expression for the joint distribution on the extended state space given in equation 18.55, we get:

$$\tilde{p}(\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{n-1}}, \mathbf{x}_1) \quad (18.57)$$

$$= \tilde{p}_1(\mathbf{x}_1) \frac{\tilde{p}_{\eta_{n-1}}(\mathbf{x}_{\eta_{n-1}})}{\tilde{p}_{\eta_{n-1}}(\mathbf{x}_1)} T_{\eta_{n-1}}(\mathbf{x}_1 | \mathbf{x}_{\eta_{n-1}}) \prod_{i=1}^{n-2} \frac{\tilde{p}_{\eta_i}(\mathbf{x}_{\eta_i})}{\tilde{p}_{\eta_i}(\mathbf{x}_{\eta_{i+1}})} T_{\eta_i}(\mathbf{x}_{\eta_{i+1}} | \mathbf{x}_{\eta_i}) \quad (18.58)$$

$$= \frac{\tilde{p}_1(\mathbf{x}_1)}{\tilde{p}_{\eta_{n-1}}(\mathbf{x}_1)} T_{\eta_{n-1}}(\mathbf{x}_1 | \mathbf{x}_{\eta_{n-1}}) \tilde{p}_{\eta_1}(\mathbf{x}_{\eta_1}) \prod_{i=1}^{n-2} \frac{\tilde{p}_{\eta_{i+1}}(\mathbf{x}_{\eta_{i+1}})}{\tilde{p}_{\eta_i}(\mathbf{x}_{\eta_{i+1}})} T_{\eta_i}(\mathbf{x}_{\eta_{i+1}} | \mathbf{x}_{\eta_i}). \quad (18.59)$$

We now have means of generating samples from the joint proposal distribution  $q$  over the extended sample via a sampling scheme given above, with the joint distribution given by

$$q(\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{n-1}}, \mathbf{x}_1) = p_0(\mathbf{x}_{\eta_1}) T_{\eta_1}(\mathbf{x}_{\eta_2} | \mathbf{x}_{\eta_1}) \dots T_{\eta_{n-1}}(\mathbf{x}_1 | \mathbf{x}_{\eta_{n-1}}). \quad (18.60)$$

We have a joint distribution on the extended space given by equation 18.59. Taking  $q(\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{n-1}}, \mathbf{x}_1)$  as the proposal distribution on the extended state space from which we will draw samples, it remains to determine the importance weights:

$$w^{(k)} = \frac{\tilde{p}(\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{n-1}}, \mathbf{x}_1)}{q(\mathbf{x}_{\eta_1}, \dots, \mathbf{x}_{\eta_{n-1}}, \mathbf{x}_1)} = \frac{\tilde{p}_1(\mathbf{x}_1^{(k)})}{\tilde{p}_{\eta_{n-1}}(\mathbf{x}_{\eta_{n-1}}^{(k)})} \dots \frac{\tilde{p}_{\eta_2}(\mathbf{x}_{\eta_2}^{(k)}) \tilde{p}_{\eta_1}(\mathbf{x}_{\eta_1}^{(k)})}{\tilde{p}_1(\mathbf{x}_{\eta_1}^{(k)}) \tilde{p}_0(\mathbf{x}_0^{(k)})}. \quad (18.61)$$

These weights are the same as proposed for AIS. Thus we can interpret AIS as simple importance sampling applied to an extended state, and its validity follows immediately from the validity of importance sampling.

Annealed importance sampling was first discovered by [Jarzynski \(1997\)](#) and then again, independently, by [Neal \(2001\)](#). It is currently the most common way of estimating the partition function for undirected probabilistic models. The reasons for this may have more to do with the publication of an influential paper ([Salakhutdinov and Murray, 2008](#)) describing its application to estimating the partition function of restricted Boltzmann machines and deep belief networks than

with any inherent advantage the method has over the other method described below.

A discussion of the properties of the AIS estimator (e.g., its variance and efficiency) can be found in [Neal \(2001\)](#).