

where  $f(\mathbf{x}, j)$  returns  $\mathbf{x}$  with the bit at position  $j$  flipped. Ratio matching avoids the partition function using the same trick as the pseudolikelihood estimator: in a ratio of two probabilities, the partition function cancels out. [Marlin \*et al.\* \(2010\)](#) found that ratio matching outperforms SML, pseudolikelihood and GSM in terms of the ability of models trained with ratio matching to denoise test set images.

Like the pseudolikelihood estimator, ratio matching requires  $n$  evaluations of  $\tilde{p}$  per data point, making its computational cost per update roughly  $n$  times higher than that of SML.

As with the pseudolikelihood estimator, ratio matching can be thought of as pushing down on all fantasy states that have only one variable different from a training example. Since ratio matching applies specifically to binary data, this means that it acts on all fantasy states within Hamming distance 1 of the data.

Ratio matching can also be useful as the basis for dealing with high-dimensional sparse data, such as word count vectors. This kind of data poses a challenge for MCMC-based methods because the data is extremely expensive to represent in dense format, yet the MCMC sampler does not yield sparse values until the model has learned to represent the sparsity in the data distribution. [Dauphin and Bengio \(2013\)](#) overcame this issue by designing an unbiased stochastic approximation to ratio matching. The approximation evaluates only a randomly selected subset of the terms of the objective, and does not require the model to generate complete fantasy samples.

See [Marlin and de Freitas \(2011\)](#) for a theoretical analysis of the asymptotic efficiency of ratio matching.

## 18.5 Denoising Score Matching

In some cases we may wish to regularize score matching, by fitting a distribution

$$p_{\text{smoothed}}(\mathbf{x}) = \int p_{\text{data}}(\mathbf{y})q(\mathbf{x} | \mathbf{y})d\mathbf{y} \quad (18.27)$$

rather than the true  $p_{\text{data}}$ . The distribution  $q(\mathbf{x} | \mathbf{y})$  is a corruption process, usually one that forms  $\mathbf{x}$  by adding a small amount of noise to  $\mathbf{y}$ .

Denoising score matching is especially useful because in practice we usually do not have access to the true  $p_{\text{data}}$  but rather only an empirical distribution defined by samples from it. Any consistent estimator will, given enough capacity, make  $p_{\text{model}}$  into a set of Dirac distributions centered on the training points. Smoothing by  $q$  helps to reduce this problem, at the loss of the asymptotic consistency property