

incompatible with bound-based positive phase methods.

### 18.3 Pseudolikelihood

Monte Carlo approximations to the partition function and its gradient directly confront the partition function. Other approaches sidestep the issue, by training the model without computing the partition function. Most of these approaches are based on the observation that it is easy to compute ratios of probabilities in an undirected probabilistic model. This is because the partition function appears in both the numerator and the denominator of the ratio and cancels out:

$$\frac{p(\mathbf{x})}{p(\mathbf{y})} = \frac{\frac{1}{Z}\tilde{p}(\mathbf{x})}{\frac{1}{Z}\tilde{p}(\mathbf{y})} = \frac{\tilde{p}(\mathbf{x})}{\tilde{p}(\mathbf{y})}. \quad (18.17)$$

The pseudolikelihood is based on the observation that conditional probabilities take this ratio-based form, and thus can be computed without knowledge of the partition function. Suppose that we partition  $\mathbf{x}$  into  $\mathbf{a}$ ,  $\mathbf{b}$  and  $\mathbf{c}$ , where  $\mathbf{a}$  contains the variables we want to find the conditional distribution over,  $\mathbf{b}$  contains the variables we want to condition on, and  $\mathbf{c}$  contains the variables that are not part of our query.

$$p(\mathbf{a} \mid \mathbf{b}) = \frac{p(\mathbf{a}, \mathbf{b})}{p(\mathbf{b})} = \frac{p(\mathbf{a}, \mathbf{b})}{\sum_{\mathbf{a}, \mathbf{c}} p(\mathbf{a}, \mathbf{b}, \mathbf{c})} = \frac{\tilde{p}(\mathbf{a}, \mathbf{b})}{\sum_{\mathbf{a}, \mathbf{c}} \tilde{p}(\mathbf{a}, \mathbf{b}, \mathbf{c})}. \quad (18.18)$$

This quantity requires marginalizing out  $\mathbf{a}$ , which can be a very efficient operation provided that  $\mathbf{a}$  and  $\mathbf{c}$  do not contain very many variables. In the extreme case,  $\mathbf{a}$  can be a single variable and  $\mathbf{c}$  can be empty, making this operation require only as many evaluations of  $\tilde{p}$  as there are values of a single random variable.

Unfortunately, in order to compute the log-likelihood, we need to marginalize out large sets of variables. If there are  $n$  variables total, we must marginalize a set of size  $n - 1$ . By the chain rule of probability,

$$\log p(\mathbf{x}) = \log p(x_1) + \log p(x_2 \mid x_1) + \cdots + \log p(x_n \mid \mathbf{x}_{1:n-1}). \quad (18.19)$$

In this case, we have made  $\mathbf{a}$  maximally small, but  $\mathbf{c}$  can be as large as  $\mathbf{x}_{2:n}$ . What if we simply move  $\mathbf{c}$  into  $\mathbf{b}$  to reduce the computational cost? This yields the **pseudolikelihood** (Besag, 1975) objective function, based on predicting the value of feature  $x_i$  given all of the other features  $\mathbf{x}_{-i}$ :

$$\sum_{i=1}^n \log p(x_i \mid \mathbf{x}_{-i}). \quad (18.20)$$