

**Algorithm 19.1 Computing the gradient in a network with table-CPDs**


---

```

Procedure Compute-Gradient (
   $\mathcal{G}$ , // Bayesian network structure over  $X_1, \dots, X_n$ 
   $\theta$ , // Set of parameters for  $\mathcal{G}$ 
   $\mathcal{D}$  // Partially observed data set
)
1 // Initialize data structures
2 for each  $i = 1, \dots, n$ 
3   for each  $x_i, \mathbf{u}_i \in \text{Val}(X_i, \text{Pa}_{X_i}^{\mathcal{G}})$ 
4      $\bar{M}[x_i, \mathbf{u}_i] \leftarrow 0$ 
5   // Collect probabilities from all instances
6 for each  $m = 1 \dots M$ 
7   Run clique tree calibration on  $\langle \mathcal{G}, \theta \rangle$  using evidence  $\mathbf{o}[m]$ 
8   for each  $i = 1, \dots, n$ 
9     for each  $x_i, \mathbf{u}_i \in \text{Val}(X_i, \text{Pa}_{X_i}^{\mathcal{G}})$ 
10       $\bar{M}[x_i, \mathbf{u}_i] \leftarrow \bar{M}[x_i, \mathbf{u}_i] + P(x_i, \mathbf{u}_i \mid \mathbf{o}[m])$ 
11   // Compute components of the gradient vector
12 for each  $i = 1, \dots, n$ 
13   for each  $x_i, \mathbf{u}_i \in \text{Val}(X_i, \text{Pa}_{X_i}^{\mathcal{G}})$ 
14      $\delta_{x_i|\mathbf{u}_i} \leftarrow \frac{1}{\theta_{x_i|\mathbf{u}_i}} \bar{M}[x_i, \mathbf{u}_i]$ 
15 return  $\{\delta_{x_i|\mathbf{u}_i} : \forall i = 1, \dots, n, \forall (x_i, \mathbf{u}_i) \in \text{Val}(X_i, \text{Pa}_{X_i}^{\mathcal{G}})\}$ 

```

---

**19.2.1.3 Gradient Ascent Algorithm**

We now generalize these ideas to case of an arbitrary network. For now we focus on the case of table-CPDs. In this case, the gradient is given by theorem 19.2. To compute the gradient for the CPD  $P(X \mid \mathbf{U})$ , we need to compute the joint probability of  $x$  and  $\mathbf{u}$  relative to our current parameter setting  $\theta$  and each observed instance  $\mathbf{o}[m]$ . In other words, we need to compute the joint distribution  $P(X[m], \mathbf{U}[m] \mid \mathbf{o}[m], \theta)$  for each  $m$ . We can do this by running an inference procedure for each data case. Importantly, **we can do all of the required inference for each data case using one clique tree calibration, since the family preservation property guarantees that  $X$  and its parents  $\mathbf{U}$  will be together in some clique in the tree.** Procedure Compute-Gradient, shown in algorithm 19.1, performs these computations.



Once we have a procedure for computing the gradient, it seems that we can simply plug it into a standard package for gradient ascent and optimize the parameters. As we have illustrated, however, there is one issue that we need to deal with. It is not hard to confirm that all components of the gradient vector are nonnegative. This is natural, since increasing each of the parameters will lead to higher likelihood. Thus, a step in the gradient direction will *increase* all the parameters. Remember, however, that we want to ensure that our parameters describe a legal probability distribution. That is, the parameters for each conditional probability are nonnegative and sum to one.