

译并解释以下内容，翻译不能有遗漏，解释请通俗易懂且详尽，因为我的数学基础不太好，最后总结要点。

## 20 Exact inference for graphical models

### 20.1 Introduction

In Section 17.4.3, we discussed the forwards-backwards algorithm, which can exactly compute the posterior marginals  $p(x_t|v, \theta)$  in any chain-structured graphical model, where  $x$  are the hidden variables (assumed discrete) and  $v$  are the visible variables. This algorithm can be modified to compute the posterior mode and posterior samples. A similar algorithm for linear-Gaussian chains, known as the Kalman smoother, was discussed in Section 18.3.2. Our goal in this chapter is to generalize these exact inference algorithms to arbitrary graphs. The resulting methods apply to both directed and undirected graphical models. We will describe a variety of algorithms, but we omit their derivations for brevity. See e.g., (Darwiche 2009; Koller and Friedman 2009) for a detailed exposition of exact inference techniques for discrete directed graphical models.

### 翻译：

## 20 精确推断图模型

### 20.1 引言

在第17.4.3节中，我们讨论了前向后向算法，该算法能够精确计算任何链式结构图模型中的后验边缘分布  $p(x_t|v, \theta)$ ，其中  $x$  是隐藏变量（假设是离散的），而  $v$  是可见变量。这个算法可以修改为计算后验众数和后验样本。类似的算法也可以用于线性高斯链，称为卡尔曼平滑器，已经在第18.3.2节中讨论过。本章的目标是将这些精确推断算法推广到任意图形模型。得到的方法适用于有向和无向图模型。我们将描述多种算法，但为了简洁起见，省略了它们的推导过程。有关离散有向图模型的精确推断技术的详细阐述，请参见 (Darwiche 2009; Koller 和 Friedman 2009)。

### 解释：

#### 1. 精确推断 (Exact Inference) :

在图模型中，推断是指从观察到的可见变量推断隐藏变量的过程。精确推断则意味着能够在数学上得出准确的结果，而不是通过近似或估计。

#### 2. 图模型 (Graphical Models) :

图模型是用于表示复杂变量之间依赖关系的结构化模型，通常包括节点（变量）和边（变量之间的依赖）。图可以是有向的，也可以是无向的。

- **有向图**表示变量之间的条件依赖关系。
- **无向图**则表示变量之间的无条件依赖关系。

### 3. 链式结构图模型：

链式结构图模型是一种特定类型的图模型，具有线性排列的结构。例如，隐马尔可夫模型（HMM）就是一个典型的链式结构图模型，其中变量沿时间轴排列，每个变量只与前一个变量和当前变量相连。

### 4. 前向后向算法（Forwards-Backwards Algorithm）：

这是一个经典的用于链式图模型的算法。它能精确地计算每个隐藏变量（例如在时间序列中）给定观察到的可见变量之后的后验分布。具体来说，它可以计算出某个时间点隐藏变量的概率。

### 5. 后验边缘分布 $p(x_t | v, \theta)$ ：

这表示在给定观察数据  $v$  和模型参数  $\theta$  的情况下，某个时刻  $t$  隐藏变量  $x_t$  的分布。

### 6. 卡尔曼平滑器（Kalman Smoother）：

这是一种特别针对线性高斯链（线性动态系统中，变量服从高斯分布）的一种推断方法，用于通过观察历史数据来估计系统的状态。

### 7. 目标：推广到任意图模型：

该章节的目的是将上述算法（例如前向后向算法和卡尔曼平滑器）推广到更一般的图模型中，无论这些图是有向图还是无向图。通过这种方式，可以应用这些推断技术解决更广泛的问题。

### 8. 省略推导：

在本节中，作者将直接给出各种算法的描述，而不详细推导每个算法的数学背景。为了理解这些算法的具体应用，作者推荐参考其他详细文献（如Darwiche 2009年、Koller和Friedman 2009年出版的书籍）。

---

## 总结要点：

1. 本章讨论如何将前向后向算法和卡尔曼平滑器推广到更一般的图模型中，适用于有向和无向图。
2. 精确推断用于通过观察到的数据推断隐藏变量。
3. 本章的目标是给出各种算法的描述，但推导过程会省略，推荐进一步阅读相关文献以获取详细的推导过程。



翻译并解释以下内容，翻译不能有遗漏，解释请通俗易懂且详尽，因为我的数学基础不太好，最后总结要点。

Figure 20.1 Message passing on a tree. (a) Collect-to-root phase. (b) Distribute-from-root phase.

Figure 20.2 Illustration of how to compute the two-slice distribution for an HMM. The  $\psi_t$  and  $\psi_{t+1}$  terms are the local evidence messages from the visible nodes  $v_t, v_{t+1}$  to the hidden nodes  $x_t, x_{t+1}$  respectively;  $f_t$  is the forwards message from  $x_{t-1}$  and  $\beta_{t+1}$  is the backwards message from  $x_{t+2}$ .

## 20.2 Belief propagation for trees

In this section, we generalize the forwards-backwards algorithm from chains to trees. The resulting algorithm is known as belief propagation (BP) (Pearl 1988), or the sum-product algorithm.

### 20.2.1 Serial protocol

We initially assume (for notational simplicity) that the model is a pairwise MRF (or CRF), i.e.,

$$p(x|v) = \frac{1}{Z(v)} \prod_{s \in V} \psi_s(x_s) \prod_{(s,t) \in E} \psi_{s,t}(x_s, x_t) \quad (20.1)$$

where  $\psi_s$  is the local evidence for node  $s$ , and  $\psi_{s,t}$  is the potential for edge  $s - t$ . We will consider the case of models with higher order cliques (such as directed trees) later on. One way to implement BP for undirected trees is as follows. Pick an arbitrary node and call it the root,  $r$ . Now orient all edges away from  $r$  (intuitively, we can imagine “picking up the graph” at node  $r$  and letting all the edges “dangle” down). This gives us a well-defined notion of parent and child. Now we send messages up from the leaves to the root (the collect evidence phase) and then back down from the root (the distribute evidence phase), in a manner analogous to forwards-backwards on chains.

To explain the process in more detail, consider the example in Figure 20.1. Suppose we want to compute the belief state at node  $t$ . We will initially condition the belief only on evidence that is at or below  $t$  in the graph, i.e., we want to compute  $\text{bel- } t(x_t) = p(x_t | v - t)$ . We will call this a “bottom-up belief state”. Suppose, by induction, that we have computed “messages” from  $t$ ’s two children, summarizing what they think  $t$  should know about the evidence in their subtrees, i.e., we have computed  $m-$

$s \rightarrow t(xt) = p(xt|v - st)$ , where  $v - st$  is all the evidence on the downstream side of the  $s - t$  edge (see Figure 20.1(a)), and similarly we have computed  $m_{+} \rightarrow t(xt)$ . Then we can compute the bottom-up belief state at  $t$  as follows:

$$\text{bel- } t(xt) = p(xt|v - t) = \frac{1}{Z_t} \prod_{c \in ch(t)} m_{-} c \rightarrow t(xt) \quad (20.2)$$

where  $\psi_t(xt) \propto p(xt|vt)$  is the local evidence for node  $t$ , and  $Z_t$  is the local normalization constant. In words, we multiply all the incoming messages from our children, as well as the incoming message from our local evidence, and then normalize. We have explained how to compute the bottom-up belief states from the bottom-up messages. How do we compute the messages themselves? Consider computing  $m_{-} s \rightarrow t(xt)$ , where  $s$  is one of  $t$ 's children. Assume, by recursion, that we have computed  $\text{bel- } s(xs) = p(xs|v - s)$ . Then we can compute the message as follows:

$$m_{-} s \rightarrow t(xt) = xs \psi_{st}(xs, xt) \text{bel- } s(xs) \quad (20.3)$$

Essentially we convert beliefs about  $xs$  into beliefs about  $xt$  by using the edge potential  $\psi_{st}$ . We continue in this way up the tree until we reach the root. Once at the root, we have "seen" all the evidence in the tree, so we can compute our local belief state at the root using  
 $\text{belr}(xr) = p(xr|v) = p(xt|v - r) \propto \psi_r(xr) \prod_{c \in ch(r)} m_{-} c \rightarrow r(xr)$

(20.4) This completes the end of the upwards pass, which is analogous to the forwards pass in an HMM. As a "side effect", we can compute the probability of the evidence by collecting the normalization constants:

$$p(v) = \prod_t Z_t \quad (20.5)$$

We can now pass messages down from the root. For example, consider node  $s$ , with parent  $t$ , as shown in Figure 20.1(b). To compute the belief state for  $s$ , we need to combine the bottom-up belief for  $s$  together with a top-down message from  $t$ , which summarizes all the information in the rest of the graph,  $m_{+} t \rightarrow s(xs) = p(xt|v + st)$ , where  $v + st$  is all the evidence on the upstream (root) side of the  $s - t$  edge, as shown in Figure 20.1(b). We then have

$$\text{bel- } s(xs) = p(xs|v) \propto \text{bel- } s(xs) \prod_{t \in pa(s)} m_{+} t \rightarrow s(xt) \quad (20.6)$$

How do we compute these downward messages? For example, consider the message from  $t$  to  $s$ . Suppose  $t$ 's parent is  $r$ , and  $t$ 's children are  $s$  and  $u$ , as shown in Figure 20.1(b). We want to include in  $m_{+} t \rightarrow s$  all the information that  $t$  has received, except for the information that  $s$  sent it:

$$m_{+} t \rightarrow s(xs) = p(xs|v + st) = \prod_{x \in ch(t)} \psi_{st}(xs, xt) \text{bel- } t(xt) m_{-} s \rightarrow t(xt) \quad (20.7)$$

Rather than dividing out the message sent up to  $t$ , we can plug in the equation of belt to get

$$m + t \rightarrow s(xs) = xt \psi_{st}(xs, xt) \psi_t(xt) c \in ch(t), c = s m - c \rightarrow t(xt) p \in pa(t) m + p \rightarrow t(xt) \quad (20.8)$$

In other words, we multiply together all the messages coming into  $t$  from all nodes except for the recipient  $s$ , combine together, and then pass through the edge potential  $\psi_{st}$ . In the case of a chain,  $t$  only has one child  $s$  and one parent  $p$ , so the above simplifies to

$$m + t \rightarrow s(xs) = xt \psi_{st}(xs, xt) \psi_t(xt) m + p \rightarrow t(xt) \quad (20.9)$$

The version of BP in which we use division is called belief updating, and the version in which we multiply all-but-one of the messages is called sum-product. The belief updating version is analogous to how we formulated the Kalman smoother in Section 18.3.2: the topdown messages depend on the bottom-up messages. This means they can be interpreted as conditional posterior probabilities. The sum-product version is analogous to how we formulated the backwards algorithm in Section 17.4.3: the top-down messages are completely independent of the bottom-up messages, which means they can only be interpreted as conditional likelihoods. See Section 18.3.2.3 for a more detailed discussion of this subtle difference.

### 20.2.2 Parallel protocol

So far, we have presented a serial version of the algorithm, in which we send messages up to the root and back. This is the optimal approach for a tree, and is a natural extension of forwards-backwards on chains.

However, as a prelude to handling general graphs with loops, we now consider a parallel version of BP. This gives equivalent results to the serial version but is less efficient when implemented on a serial machine.

The basic idea is that all nodes receive messages from their neighbors in parallel, they then update their belief states, and finally they send new messages back out to their neighbors. This process repeats until convergence. This kind of computing architecture is called a systolic array, due to its resemblance to a beating heart. More precisely, we initialize all messages to the all 1's vector. Then, in parallel, each node absorbs messages from all its neighbors using

$$bel_s(xs) \propto \psi_s(xs) t \in nbrs m \rightarrow s(xs) \quad (20.10)$$

Then, in parallel, each node sends messages to each of its neighbors:

$$m_s \rightarrow t(xt) = xs / (\psi_s(xs) \psi_{st}(xs, xt) u \in nbrs \setminus t mu \rightarrow s(xs)) \quad (20.11)$$

The  $m_s \rightarrow t$  message is computed by multiplying together all incoming messages, except the one sent by the recipient, and then passing through the  $\psi_{st}$  potential. At iteration  $T$  of the algorithm,  $b_{st}(x_s)$  represents the posterior belief of  $x_s$  conditioned on the evidence that is  $T$  steps away in the graph. After  $D(G)$  steps, where  $D(G)$  is the diameter of the graph (the largest distance between any two pairs of nodes), every node has obtained information from all the other nodes. Its local belief state is then the correct posterior marginal. Since the diameter of a tree is at most  $|V| - 1$ , the algorithm converges in a linear number of steps. We can actually derive the up-down version of the algorithm by imposing the condition that a node can only send a message once it has received messages from all its other neighbors. This means we must start with the leaf nodes, which only have one neighbor. The messages then propagate up to the root and back. We can also update the nodes in a random order. The only requirement is that each node get updated  $D(G)$  times. This is just enough time for information to spread throughout the whole tree. Similar parallel, distributed algorithms for solving linear systems of equations are discussed in (Bertsekas 1997). In particular, the Gauss-Seidel algorithm is analogous to the serial up-down version of BP, and the Jacobi algorithm is analogous to the parallel version of BP.

### 20.2.3 Gaussian BP

Now consider the case where  $p(x|v)$  is jointly Gaussian, so it can be represented as a Gaussian pairwise MRF, as in Section 19.4.4. We now present the belief propagation algorithm for this class of models, follow the presentation of (Bickson 2009) (see also (Malioutov et al. 2006)). We will assume the following node and edge potentials:

$$\psi_t(x_t) = \exp(-\frac{1}{2} \text{Att}x_t^2 t + b_t x_t) \quad (20.12)$$

$$\psi_{st}(x_s, x_t) = \exp(-\frac{1}{2} x_s A s t x_t) \quad (20.13)$$

so the overall model has the form

$$p(x|v) \propto \exp(-\frac{1}{2} x^T A x + b^T x) \quad (20.14)$$

This is the information form of the MVN (see Exercise 9.2), where  $A$  is the precision matrix. Note that by completing the square, the local evidence can be rewritten as a Gaussian:

$$\psi_t(x_t) \propto N(b_t / \text{Att}, A - 1/t) = N(m_t, I - 1/t) \quad (20.15)$$

Below we describe how to use BP to compute the posterior node marginals,

$$p(x_t|v) = N(\mu_t, \lambda - 1/t) \quad (20.16)$$

If the graph is a tree, the method is exact. If the graph is loopy, the posterior means may still be exact, but the

posterior variances are often too small (Weiss and Freeman 1999). Although the precision matrix  $A$  is often sparse, computing the posterior mean requires inverting it, since  $\mu = A^{-1}b$ . BP provides a way to exploit graph structure to perform this computation in  $O(D)$  time instead of  $O(D^3)$ . This is related to various methods from linear algebra, as discussed in (Bickson 2009). Since the model is jointly Gaussian, all marginals and all messages will be Gaussian. The key operations we need are to multiply together two Gaussian factors, and to marginalize out a variable from a joint Gaussian factor. For multiplication, we can use the fact that the product of two Gaussians is Gaussian:

$$N(x|\mu_1, \lambda_1^{-1}) \times N(x|\mu_2, \lambda_2^{-1}) = CN(x|\mu, \lambda^{-1}) \quad (20.17)$$

$$\lambda = \lambda_1 + \lambda_2 \quad (20.18)$$

$$\mu = \lambda^{-1}(\mu_1\lambda_1 + \mu_2\lambda_2) \quad (20.19)$$

where

$$C = \sqrt{\lambda_1\lambda_2} \exp(-\frac{1}{2}(\lambda_1\mu_2^2 + \lambda_2\mu_1^2 - 2\lambda_1\lambda_2\mu_1\mu_2)) \quad (20.20)$$

See Exercise 20.2 for the proof. For marginalization, we have the following result:

$$\exp(-ax^2 + bx)dx = \sqrt{\pi/a} \exp(b^2/4a) \quad (20.21)$$

which follows from the normalization constant of a Gaussian (Exercise 2.11). We now have all the pieces we need. In particular, let the message  $m_{s \rightarrow t}(x_t)$  be a Gaussian with mean  $\mu_{st}$  and precision  $\lambda_{st}$ . From Equation 20.10, the belief at node  $s$  is given by the product of incoming messages times the local evidence (Equation 20.15) and hence

$$bel_s(x_s) = \prod_{t \in nbr(s)} m_{t \rightarrow s}(x_s) \quad (20.22)$$

$$\lambda_s = \sum_{t \in nbr(s)} \lambda_{ts} \quad (20.23)$$

$$\mu_s = \lambda_s^{-1} \left( \sum_{t \in nbr(s)} \lambda_{ts} \mu_{ts} \right) \quad (20.24)$$

To compute the messages themselves, we use Equation 20.11, which is given by

$$m_{s \rightarrow t}(x_t) = \int_{-\infty}^{\infty} \psi_{st}(x_s, x_t) \psi_s(x_s) \prod_{u \in nbr(s) \setminus t} m_{u \rightarrow s}(x_u) dx_s \quad (20.25)$$

$$= \int_{-\infty}^{\infty} \psi_{st}(x_s, x_t) f_{s \setminus t}(x_s) dx_s \quad (20.26)$$

where  $f_{s \setminus t}(x_s)$  is the product of the local evidence and all incoming messages excluding the message from  $t$ :

$$f_{s \setminus t}(x_s) = \prod_{u \in nbr(s) \setminus t} m_{u \rightarrow s}(x_s) \quad (20.27)$$

$$= N(x_s | \mu_s, \lambda_s^{-1}) \quad (20.28)$$

$$\lambda_s = \sum_{u \in nbr(s) \setminus t} \lambda_{us} \quad (20.29)$$

$$\mu_s = \lambda_s^{-1} \left( \sum_{u \in nbr(s) \setminus t} \lambda_{us} \mu_{us} \right) \quad (20.30)$$

Returning to Equation 20.26 we have

$$m_{s \rightarrow t}(x_t) = \int_{-\infty}^{\infty} \exp(-\frac{1}{2}\lambda_s(x_s - \mu_s)^2 / \lambda_s) \psi_{st}(x_s, x_t) \exp(-\frac{1}{2}\lambda_t(x_t - \mu_t)^2 / \lambda_t) f_{s \setminus t}(x_s) dx_s \quad (20.31)$$

$$= \int_{-\infty}^{\infty} \exp(-\frac{1}{2}\lambda_s(x_s - \mu_s)^2 / \lambda_s) + (\lambda_s \mu_s - \lambda_t \mu_t) x_s + \frac{1}{2}(\lambda_s + \lambda_t) x_s^2 dx_s + const$$

(20.32)

$$\propto \exp(\lambda s \|\mu_s - \alpha_{st}\|^2 / (2\lambda s)) \quad (20.33)$$

$$\propto N(\mu_s, \lambda^{-1} I) \quad (20.34)$$

$$\lambda s = A^2 s / \lambda s \quad (20.35)$$

$$\mu_s = \alpha_{st} \mu_s / \lambda s \quad (20.36)$$

One can generalize these equations to the case where each node is a vector, and the messages become small MVNs instead of scalar Gaussians (Alag and Agogino 1996). If we apply the resulting algorithm to a linear dynamical system, we recover the Kalman smoothing algorithm of Section 18.3.2.

To perform message passing in models with non-Gaussian potentials, one can use sampling methods to approximate the relevant integrals. This is called non-parametric BP (Sudderth et al. 2003; Isard 2003; Sudderth et al. 2010).

#### 20.2.4 Other BP variants

In this section, we briefly discuss several variants of the main algorithm.

##### 20.2.4.1 Max-product algorithm

It is possible to devise a max-product version of the BP algorithm, by replacing the operator with the max operator. We can then compute the local MAP marginal of each node. However, if there are ties, this might not be globally consistent, as discussed in Section 17.4.4.

Fortunately, we can generalize the Viterbi algorithm to trees, where we use max and argmax in the collect-to-root phase, and perform traceback in the distribute-from-root phase. See (Dawid 1992) for details.

##### 20.2.4.2 Sampling from a tree

It is possible to draw samples from a tree structured model by generalizing the forwards filtering / backwards sampling algorithm discussed in Section 17.4.5. See (Dawid 1992) for details.

##### 20.2.4.3 Computing posteriors on sets of variables

In Section 17.4.3.2, we explained how to compute the “two-slice” distribution  $\xi_{t,t+1}(i, j) = p(x_t = i, x_{t+1} = j | v)$  in an HMM, namely by using

$$\xi_{t,t+1}(i, j) = \alpha_t(i) \psi_{t+1}(j) \beta_{t+1}(j) \psi_t(i, j) \quad (20.37)$$

Since  $\alpha_t(i) \propto \psi_t(i) f_t(i)$ , where  $f_t = p(x_t | v_{1:t-1})$  is the forwards message, we can think of this as sending messages  $f_t$  and  $\psi_t$  into  $x_t$ ,  $\beta_{t+1}$  and  $\psi_{t+1}$  into  $x_{t+1}$ , and then combining them with the  $\Psi$  matrix, as shown in Figure 20.2. This is like treating  $x_t$  and  $x_{t+1}$  as a single “mega node”, and then multiplying all the incoming messages as well as all the local factors (here,  $\psi_t, \psi_{t+1}$ ).

It looks like you've uploaded two images, but they are not directly viewable in this context. If you need me to help with specific details from the images or to provide further explanations based on them, please describe the content or let me know what you'd like me to do with these images!

翻译并解释

看起来你上传了两张图片，但由于我无法直接查看图像内容，请提供图片的描述或提问，我会尽力帮助你进行翻译和解释。

翻译并解释以下内容，翻译不能有遗漏，解释请通俗易懂且详尽，因为我的数学基础不太好，最后总结要点。

Figure 20.1 Message passing on a tree. (a) Collect-to-root phase. (b) Distribute-from-root phase.

Figure 20.2 Illustration of how to compute the two-slice distribution for an HMM. The  $\psi_t$  and  $\psi_{t+1}$  terms are the local evidence messages from the visible nodes  $v_t, v_{t+1}$  to the hidden nodes  $x_t, x_{t+1}$  respectively;  $f_t$  is the forwards message from  $x_{t-1}$  and  $\beta_{t+1}$  is the backwards message from  $x_{t+2}$ .

## 20.2 Belief propagation for trees

In this section, we generalize the forwards-backwards algorithm from chains to trees. The resulting algorithm is known as belief propagation (BP) (Pearl 1988), or the sum-product algorithm.

### 20.2.1 Serial protocol

We initially assume (for notational simplicity) that the model is a pairwise MRF (or CRF), i.e.,

$$p(x|v) = \prod_{s \in V} Z(v) \prod_{(s,t) \in E} \psi_{s,t}(x_s, x_t) \quad (20.1)$$

where  $\psi_s$  is the local evidence for node  $s$ , and  $\psi_{s,t}$  is the potential for edge  $s - t$ . We will consider the case of models with higher order cliques (such as directed trees) later on. One way to implement BP for undirected trees is as follows. Pick an arbitrary node and call it the root,  $r$ .

Now orient all edges away from  $r$  (intuitively, we can imagine “picking up the graph” at node  $r$  and letting all the edges “dangle” down). This gives us a well-defined

notion of parent and child. Now we send messages up from the leaves to the root (the collect evidence phase) and then back down from the root (the distribute evidence phase), in a manner analogous to forwards-backwards on chains.

To explain the process in more detail, consider the example in Figure 20.1. Suppose we want to compute the belief state at node  $t$ . We will initially condition the belief only on evidence that is at or below  $t$  in the graph, i.e., we want to compute  $\text{bel- } t(xt) = p(xt|v- t)$ . We will call this a "bottom-up belief state". Suppose, by induction, that we have computed "messages" from  $t$ 's two children, summarizing what they think  $t$  should know about the evidence in their subtrees, i.e., we have computed  $m- s \rightarrow t(xt) = p(xt|v- st)$ , where  $v- st$  is all the evidence on the downstream side of the  $s - t$  edge (see Figure 20.1(a)), and similarly we have computed  $m_{\mu} \rightarrow t(xt)$ . Then we can compute the bottom-up belief state at  $t$  as follows:

$$\text{bel- } t(xt) = p(xt|v- t) = 1 Z_t \psi_t(xt) c \in ch(t) m- c \rightarrow t(xt) \quad (20.2)$$

where  $\psi_t(xt) \propto p(xt|vt)$  is the local evidence for node  $t$ , and  $Z_t$  is the local normalization constant. In words, we multiply all the incoming messages from our children, as well as the incoming message from our local evidence, and then normalize. We have explained how to compute the bottom-up belief states from the bottom-up messages. How do we compute the messages themselves? Consider computing  $m- s \rightarrow t(xt)$ , where  $s$  is one of  $t$ 's children. Assume, by recursion, that we have computed  $\text{bel- } s(xs) = p(xs|v- s)$ . Then we can compute the message as follows:

$$m- s \rightarrow t(xt) = xs \psi_{st}(xs, xt) \text{bel- } s(xs) \quad (20.3)$$

Essentially we convert beliefs about  $xs$  into beliefs about  $xt$  by using the edge potential  $\psi_{st}$ . We continue in this way up the tree until we reach the root. Once at the root, we have "seen" all the evidence in the tree, so we can compute our local belief state at the root using

$$\text{bel}_r(xr) = p(xr|v) = p(xt|v- r) \propto \psi_r(xr) c \in ch(r) m- c \rightarrow r(xr) \quad (20.4)$$

This completes the end of the upwards pass, which is analogous to the forwards pass in an HMM. As a "side effect", we can compute the probability of the evidence by collecting the normalization constants:

$$p(v) = t Z_t \quad (20.5)$$

We can now pass messages down from the root. For example, consider node  $s$ , with parent  $t$ , as shown in Figure 20.1(b). To compute the belief state for  $s$ , we need to combine the bottom-up belief for  $s$  together with a

top-down message from  $t$ , which summarizes all the information in the rest of the graph,  $m+ t \rightarrow s(xs) = p(xt|v+ st)$ , where  $v+ st$  is all the evidence on the upstream (root) side of the  $s - t$  edge, as shown in Figure 20.1(b). We then have

$$bel_s(xs) = p(xs|v) \propto bel_{-s}(xs) \quad t \in pa(s) \quad m+ t \rightarrow s(xs) \quad (20.6)$$

How do we compute these downward messages? For example, consider the message from  $t$  to  $s$ . Suppose  $t$ 's parent is  $r$ , and  $t$ 's children are  $s$  and  $u$ , as shown in Figure 20.1(b). We want to include in  $m+ t \rightarrow s$  all the information that  $t$  has received, except for the information that  $s$  sent it:

$$m+ t \rightarrow s(xs) = p(xs|v+ st) = xt \psi_{st}(xs, xt) belt(xt) m- s \rightarrow t(xs) \quad (20.7)$$

Rather than dividing out the message sent up to  $t$ , we can plug in the equation of  $belt$  to get

$$m+ t \rightarrow s(xs) = xt \psi_{st}(xs, xt) \psi_t(xt) c \in ch(t), c = s \quad m- c \rightarrow t(xs) \quad p \in pa(t) \quad m+ p \rightarrow t(xs) \quad (20.8)$$

In other words, we multiply together all the messages coming into  $t$  from all nodes except for the recipient  $s$ , combine together, and then pass through the edge potential  $\psi_{st}$ . In the case of a chain,  $t$  only has one child  $s$  and one parent  $p$ , so the above simplifies to

$$m+ t \rightarrow s(xs) = xt \psi_{st}(xs, xt) \psi_t(xt) m+ p \rightarrow t(xs) \quad (20.9)$$

The version of BP in which we use division is called belief updating, and the version in which we multiply all-but-one of the messages is called sum-product. The belief updating version is analogous to how we formulated the Kalman smoother in Section 18.3.2: the topdown messages depend on the bottom-up messages. This means they can be interpreted as conditional posterior probabilities. The sum-product version is analogous to how we formulated the backwards algorithm in Section 17.4.3: the top-down messages are completely independent of the bottom-up messages, which means they can only be interpreted as conditional likelihoods. See Section 18.3.2.3 for a more detailed discussion of this subtle difference.

## 20.2.2 Parallel protocol

So far, we have presented a serial version of the algorithm, in which we send messages up to the root and back. This is the optimal approach for a tree, and is a natural extension of forwards-backwards on chains. However, as a prelude to handling general graphs with loops, we now consider a parallel version of BP. This gives equivalent results to the serial version but is less efficient when implemented on a serial machine.

The basic idea is that all nodes receive messages from their neighbors in parallel, they then update their belief states, and finally they send new messages back out to their neighbors. This process repeats until convergence. This kind of computing architecture is called a systolic array, due to its resemblance to a beating heart. More precisely, we initialize all messages to the all 1's vector. Then, in parallel, each node absorbs messages from all its neighbors using

$$\text{bels}(xs) \propto \psi_s(xs) \ t \in \text{nbrs} \ \text{mt} \rightarrow s(xs) \quad (20.10)$$

Then, in parallel, each node sends messages to each of its neighbors:

$$ms \rightarrow t(xt) = xs / (\psi_s(xs) \psi_{st}(xs, xt) \ u \in \text{nbrs} \setminus t \ \text{mu} \rightarrow s(xs)) \quad (20.11)$$

The  $ms \rightarrow t$  message is computed by multiplying together all incoming messages, except the one sent by the recipient, and then passing through the  $\psi_{st}$  potential. At iteration T of the algorithm,  $\text{bels}(xs)$  represents the posterior belief of  $xs$  conditioned on the evidence that is T steps away in the graph. After  $D(G)$  steps, where  $D(G)$  is the diameter of the graph (the largest distance between any two pairs of nodes), every node has obtained information from all the other nodes. Its local belief state is then the correct posterior marginal. Since the diameter of a tree is at most  $|V| - 1$ , the algorithm converges in a linear number of steps. We can actually derive the up-down version of the algorithm by imposing the condition that a node can only send a message once it has received messages from all its other neighbors. This means we must start with the leaf nodes, which only have one neighbor. The messages then propagate up to the root and back. We can also update the nodes in a random order. The only requirement is that each node get updated  $D(G)$  times. This is just enough time for information to spread throughout the whole tree. Similar parallel, distributed algorithms for solving linear systems of equations are discussed in (Bertsekas 1997). In particular, the Gauss-Seidel algorithm is analogous to the serial up-down version of BP, and the Jacobi algorithm is analogous to the parallel version of BP.

### 20.2.3 Gaussian BP

Now consider the case where  $p(x|v)$  is jointly Gaussian, so it can be represented as a Gaussian pairwise MRF, as in Section 19.4.4. We now present the belief propagation algorithm for this class of models, follow the presentation of (Bickson 2009) (see also (Malioutov et al. 2006)). We will assume the following node and edge potentials:

$$\psi_t(xt) = \exp(-\frac{1}{2} \text{Att}x^2 t + bxt) \quad (20.12)$$

$$\psi_{st}(xs, xt) = \exp(-\frac{1}{2} xsAxt) \quad (20.13)$$

so the overall model has the form

$$p(x|v) \propto \exp(-\frac{1}{2} xTAx + bTx) \quad (20.14)$$

This is the information form of the MVN (see Exercise 9.2), where A is the precision matrix. Note that by completing the square, the local evidence can be rewritten as a Gaussian:

$$\psi_t(xt) \propto N(bt/\text{Att}, A - 1/t) = N(mt, I - 1/t) \quad (20.15)$$

Below we describe how to use BP to compute the posterior node marginals,

$$p(xt|v) = N(\mu_t, \lambda - 1/t) \quad (20.16)$$

If the graph is a tree, the method is exact. If the graph is loopy, the posterior means may still be exact, but the posterior variances are often too small (Weiss and Freeman 1999). Although the precision matrix A is often sparse, computing the posterior mean requires inverting it, since  $\mu = A^{-1}b$ . BP provides a way to exploit graph structure to perform this computation in  $O(D)$  time instead of  $O(D^3)$ . This is related to various methods from linear algebra, as discussed in (Bickson 2009). Since the model is jointly Gaussian, all marginals and all messages will be Gaussian. The key operations we need are to multiply together two Gaussian factors, and to marginalize out a variable from a joint Gaussian factor. For multiplication, we can use the fact that the product of two Gaussians is Gaussian:

$$N(x|\mu_1, \lambda_1) \times N(x|\mu_2, \lambda_2) = CN(x|\mu, \lambda) \quad (20.17)$$

$$\lambda = \lambda_1 + \lambda_2 \quad (20.18)$$

$$\mu = \lambda^{-1}(\mu_1\lambda_1 + \mu_2\lambda_2) \quad (20.19)$$

where

$$C = \lambda_1\lambda_2 \exp(1/2(\lambda_1\mu_2^2/(\lambda_1 + \lambda_2) - \lambda_2\mu_1^2/(\lambda_1 + \lambda_2))) + \lambda_2\mu_2^2/(\lambda_1 + \lambda_2) \quad (20.20)$$

See Exercise 20.2 for the proof. For marginalization, we have the following result:

$$\exp(-ax^2 + bx)dx = \pi/a \exp(b^2/4a) \quad (20.21)$$

which follows from the normalization constant of a Gaussian (Exercise 2.11). We now have all the pieces we need. In particular, let the message  $m_{st} \rightarrow t(xt)$  be a Gaussian with mean  $\mu_{st}$  and precision  $\lambda_{st}$ . From Equation 20.10, the belief at node s is given by the product of incoming messages times the local evidence (Equation 20.15) and hence

$$bel_s(xs) = \psi_s(xs) \prod_{t \in nbr(s)} m_{ts}(xs) = N(xs|\mu_s, \lambda_s) \quad (20.22)$$

$$\lambda_s = s + \sum_{t \in nbr(s)} \lambda_{ts} \quad (20.23)$$

$$\mu_s = \lambda_s^{-1} \left( \sum_{t \in nbr(s)} \lambda_{ts} \mu_{ts} \right) \quad (20.24)$$

To compute the messages themselves, we use Equation

20.11, which is given by

$$ms \rightarrow t(xt) = xs / (\psi_{st}(xs, xt) \psi_s(xs) \cup_{u \in nbrs \setminus t} mu \rightarrow s(xs)) \int dx_s \quad (20.25)$$

$$= xs \psi_{st}(xs, xt) fs \setminus t(xs) dx_s \quad (20.26)$$

where  $fs \setminus t(xs)$  is the product of the local evidence and all incoming messages excluding the message from  $t$ :

$$fs \setminus t(xs) = \psi_s(xs) \cup_{u \in nbrs \setminus t} mu \rightarrow s(xs) \quad (20.27)$$

$$= N(xs | \mu_s(t, \lambda - 1) s(t)) \quad (20.28)$$

$$\lambda s \setminus t = s + u \in nbr(s) \setminus t \lambda u s \quad (20.29)$$

$$\mu s \setminus t = \lambda - 1 s(t) / ( \sum_{u \in nbr(s) \setminus t} \lambda u s \mu u ) \quad (20.30)$$

Returning to Equation 20.26 we have

$$ms \rightarrow t(xt) = xs \exp(-xs A s t xt) / \prod_{u \in nbrs \setminus t} \psi_{st}(xs, xt) \exp(-\lambda s \setminus t / 2 (xs - \mu s \setminus t)^2) / \prod_{u \in nbrs \setminus t} fs \setminus t(xs) dx_s \quad (20.31)$$

$$= xs \exp(-\lambda s \setminus t x^2 / 2) + (\lambda s \setminus t \mu s \setminus t - A s t xt) xs dx_s + \text{const} \quad (20.32)$$

$$\propto \exp(\lambda s \setminus t \mu s \setminus t - A s t xt)^2 / (2\lambda s \setminus t) \quad (20.33)$$

$$\propto N(\mu_{st}, \lambda - 1 s(t)) \quad (20.34)$$

$$\lambda s t = A s t / \lambda s \setminus t \quad (20.35)$$

$$\mu_{st} = A s t \mu s \setminus t / \lambda s t \quad (20.36)$$

One can generalize these equations to the case where each node is a vector, and the messages become small MVNs instead of scalar Gaussians (Alag and Agogino 1996). If we apply the resulting algorithm to a linear dynamical system, we recover the Kalman smoothing algorithm of Section 18.3.2.

To perform message passing in models with non-Gaussian potentials, one can use sampling methods to approximate the relevant integrals. This is called non-parametric BP (Sudderth et al. 2003; Isard 2003; Sudderth et al. 2010).

#### 20.2.4 Other BP variants

In this section, we briefly discuss several variants of the main algorithm.

##### 20.2.4.1 Max-product algorithm

It is possible to devise a max-product version of the BP algorithm, by replacing the operator with the max operator. We can then compute the local MAP marginal of each node. However, if there are ties, this might not be globally consistent, as discussed in Section 17.4.4.

Fortunately, we can generalize the Viterbi algorithm to trees, where we use max and argmax in the collect-to-root phase, and perform traceback in the distribute-from-root phase. See (Dawid 1992) for details.

##### 20.2.4.2 Sampling from a tree

It is possible to draw samples from a tree structured model by generalizing the forwards filtering / backwards sampling algorithm discussed in Section 17.4.5. See (Dawid 1992) for details.

### 20.2.4.3 Computing posteriors on sets of variables

In Section 17.4.3.2, we explained how to compute the “two-slice” distribution  $\xi_{t,t+1}(i,j) = p(x_t = i, x_{t+1} = j | v)$  in an HMM, namely by using

$$\xi_{t,t+1}(i,j) = \alpha_t(i)\psi_{t+1}(j)\beta_{t+1}(j)\psi_{t,t+1}(i,j) \quad (20.37)$$

Since  $\alpha_t(i) \propto \psi_t(i)f_t(i)$ , where  $f_t = p(x_t | v_1:t-1)$  is the forwards message, we can think of this as sending messages  $f_t$  and  $\psi_t$  into  $x_t$ ,  $\beta_{t+1}$  and  $\psi_{t+1}$  into  $x_{t+1}$ , and then combining them with the  $\Psi$  matrix, as shown in Figure 20.2. This is like treating  $x_t$  and  $x_{t+1}$  as a single “mega node”, and then multiplying all the incoming messages as well as all the local factors (here,  $\psi_{t,t+1}$ ).

## 翻译与解释：

### 图20.1 树上的信息传递

- (a) 从叶子到根的收集阶段。
- (b) 从根到叶子的分发阶段。

### 图20.2 计算HMM的两切片分布的示意图

其中， $\psi_t$  和  $\psi_{t+1}$  是从可见节点  $v_t$  和  $v_{t+1}$  到隐藏节点  $x_t$  和  $x_{t+1}$  的局部证据消息； $f_t$  是从  $x_{t-1}$  传来的前向消息， $\beta_{t+1}$  是从  $x_{t+2}$  传来的后向消息。

## 20.2 树的信念传播

本节内容是将链式结构上的前向后向算法推广到树结构。所得到的算法被称为信念传播 (Belief Propagation, BP)，也叫做和积算法 (Sum-Product Algorithm)。

### 20.2.1 串行协议

我们最初假设（为了简化符号）模型是一个成对的马尔可夫随机场 (MRF) 或条件随机场 (CRF)，即：

$$p(x|v) = \frac{1}{Z(v)} \prod_{s \in V} \psi_s(x_s) \prod_{(s,t) \in E} \psi_{s,t}(x_s, x_t)$$

其中， $\psi_s$  是节点  $s$  的局部证据， $\psi_{st}$  是边  $s$  到  $t$  的潜力。

一种实现无向树的信念传播 (BP) 的方法是：选择一个任意节点作为根节点  $r$ ，并将所有边的方向指向根节点（可以想象是“把图提起来”，让所有边垂直向下）。这就给我们一个明确的父子关系。

接下来，我们将消息从叶节点发送到根节点（收集证据阶段），然后从根节点再发送回叶节点（分发证据阶段），这与链上的前向后向过程类似。

以图20.1为例，假设我们要计算节点  $t$  的信念状态。我们最初只根据图中在  $t$  以下或等于  $t$  的证据来计算信念状态，即我们要计算：

$$\text{bel}_t^-(x_t) = p(x_t | v_t^-)$$

这个过程被称为“自下而上的信念状态”。假设我们通过递归计算了从  $t$  的两个子节点发送的消息，分别总结了它们认为  $t$  应该知道的子树证据，即我们已经计算了：

$$m_{s \rightarrow t}^-(x_t) = p(x_t | v_{st}^-)$$

其中， $v_{st}^-$  是边  $s$  到  $t$  的下游所有证据（见图20.1(a)）。然后我们可以根据以下公式计算节点  $t$  的自下而上的信念状态：

$$\text{bel}_t^-(x_t) = p(x_t | v_t^-) = \frac{1}{Z_t} \psi_t(x_t) \prod_{c \in \text{ch}(t)} m_{c \rightarrow t}^-(x_t)$$

这里， $\psi_t(x_t) \propto p(x_t | v_t)$  是节点  $t$  的局部证据， $Z_t$  是本地的归一化常数。简而言之，我们将所有来自子节点的消息与本地证据的消息相乘，并进行归一化。

---

## 消息计算

如何计算这些消息？例如，考虑计算从子节点  $s$  到节点  $t$  的消息  $m_{s \rightarrow t}^-(x_t)$ 。假设我们通过递归计算了子节点  $s$  的信念状态：

$$\text{bel}_s^-(x_s) = p(x_s | v_s^-)$$

然后消息计算如下：

$$m_{s \rightarrow t}^-(x_t) = \sum_{x_s} \psi_{st}(x_s, x_t) \text{bel}_s^-(x_s)$$

本质上，我们通过边潜力  $\psi_{st}$  将节点  $s$  的信念转化为节点  $t$  的信念。

当我们向树的根节点逐步传递消息后，到达根节点时，我们就能够“看到”树中的所有证据，从而计算根节点的本地信念状态：

$$\text{bel}_r(x_r) = p(x_r | v) = p(x_r | v_r^-) \propto \psi_r(x_r) \prod_{c \in \text{ch}(r)} m_{c \rightarrow r}^-(x_r)$$

此时，我们已经完成了“向上传递”的过程，它类似于隐马尔可夫模型中的前向过程。作为副产品，我们可以通过收集归一化常数来计算证据的概率：

$$p(v) = \prod_t Z_t$$


---

## 消息从根节点向下传递

接下来，考虑从根节点向下传递消息。以节点  $s$  为例，假设它的父节点是  $t$ 。我们需要将从根节点传来的消息与“自下而上”的信念状态结合，来计算节点  $s$  的信念状态。

消息的计算方式如下：

$$m_{t \rightarrow s}^+(x_s) = \sum_{x_t} \psi_{st}(x_s, x_t) \text{bel}_t(x_t) m_{s \rightarrow t}^-(x_t)$$


---

### 20.2.2 并行协议

到目前为止，我们介绍了一个串行的消息传递算法，其中消息先向根节点传递，再向叶节点传递。这种方式对于树结构是最优的，但对于具有循环的图来说，我们需要引入并行版本的信念传播算法。

并行协议的基本思想是：所有节点都同时接收来自邻居的消息，更新信念状态，然后将新的消息发送回邻居。这个过程会重复进行，直到收敛。

---

### 20.2.3 高斯BP

如果  $p(x|v)$  是联合高斯分布，可以通过高斯成对的马尔可夫随机场（MRF）来表示，并且可以使用信念传播来计算后验节点边缘分布。

当图是树时，该方法是精确的。如果图是带环的，后验均值可能依然是精确的，但后验方差通常太小。

---

### 20.2.4 其他BP变体

1. **最大积算法 (Max-product algorithm)**：可以通过将操作符替换为最大值操作符来设计最大积版本的信念传播算法，从而计算每个节点的局部最大后验 (MAP) 边缘。
  2. **从树中采样 (Sampling from a tree)**：通过推广前向过滤/后向采样算法，可以从树结构模型中采样。
  3. **计算变量集的后验**：可以通过类似的方式计算HMM的两切片分布，结合前向消息和后向消息来进行。
- 

### 总结要点：

1. **信念传播 (BP)** 是一种推断算法，用于计算图模型中各个节点的后验分布。它通过在树结构上传递消息来进行推断。
2. **信念状态** 是通过从子节点到父节点的“自下而上”阶段和从父节点到子节点的“自上而下”阶段相结合计算的。

3. **并行协议** 使得在图的每个节点能够并行地接收和传递消息，直到达到收敛状态，这为处理带有环的图提供了一个基础。
4. **高斯BP** 针对联合高斯分布进行了优化，使得可以更高效地计算节点的后验均值和方差。

这一节主要在讲什么？我不太懂

这节主要在讲\*\*信念传播（Belief Propagation, BP）\*\*算法，它是一种用于计算图模型中各个节点后验概率的算法。简单来说，它是通过在图的各个节点之间传递消息来推断隐藏的变量或状态。

## 信念传播的背景

信念传播最初是为了解决在树结构（没有循环的图）上推断的问题。假设我们有一个复杂的图，节点代表不同的变量（比如某种状态或信息），边代表这些变量之间的关系。BP算法的目标是帮助我们从已知的信息（可见变量）推断出未知的信息（隐藏变量）。

## 树结构与消息传递

- **树结构**: 树是一个没有循环的图，其中有一个“根”节点，其他节点通过边连接。树结构的特点是，消息从一个节点传递到另一个节点时，路径上没有回路。
- **消息传递**: 在信念传播中，消息是指一个节点通过其边与相邻节点交换的信念。每个节点会根据邻居节点传来的消息和自己的证据（本地信息）来更新自己的信念。这些信念就可以用来推测节点的状态或变量值。

## 这节内容的核心：

### 1. 信念传播的工作原理：

- 在一个树形结构中，我们从叶节点（树的末端节点）开始传递消息，直到根节点。然后，根节点开始将信息反向传递给叶节点。
- 消息传递的过程分为两部分：
  - **收集证据阶段**: 从叶子节点向根节点传递信息。
  - **分发证据阶段**: 从根节点向叶子节点传递信息。

### 2. 如何计算信念：

- 信念传播通过计算各个节点的“信念状态”来完成推断。信念状态就是节点在给定邻居信息后的“最可能的状态”。
- 每个节点在计算自己的信念时，会考虑从它的子节点或父节点传来的信息。通过逐步更新信念，最终会得到图中每个节点的后验分布。

### 3. 并行与串行版本：

- **串行版本**: 我们按顺序传递消息，从叶节点开始到根节点，再从根节点回传到叶节点。

- **并行版本**: 所有节点同时接收并传递消息，这对于处理更复杂的图（带有循环的图）更为高效。

#### 4. 高斯信念传播:

- 对于高斯分布的情况，信念传播可以用来计算每个节点的后验均值和方差，特别是在高斯模型中非常有效。

#### 总结:

- 这节的核心在于**信念传播**，通过在图的节点之间传递消息来计算节点的信念状态，从而推断出节点的隐藏信息。信念传播可以处理树结构图，也可以扩展到带环的图，并且在高斯模型中具有特殊的优化方法。

希望这个解释能帮你理解这节的核心内容！如果你对其中的某个概念还有疑问，或者需要更详细的说明，请告诉我。

请先翻译再解释以下内容，翻译不能有遗漏，解释请通俗易懂且详尽，因为我的数学基础不太好，最后总结要点。  
强调：翻译不能有遗漏，必须全部翻译。

Figure 20.3 Left: The “student” DGM. Right: the equivalent UGM. We add moralization arcs D-I, G-J and L-S.

Figure 20.4 Example of the elimination process, in the order C, D, I, etc. When we eliminate I (figure c), we add a fill-in edge between G and S, since they are not connected.

Figure 20.5 Sending multiple messages along a tree. (a) X1 is root. (b) X2 is root. (c) X4 is root. (d) All of the messages needed to compute all singleton marginals.

Table 20.1 Eliminating variables from Figure 20.3 in the order C,D, I,H, G, S, L to compute P(J).

Table 20.2 Some commutative semirings.

Table 20.3 Eliminating variables from Figure 20.3 in the order G, I, S, L,H, C,D.

20.3 The variable elimination algorithm

We have seen how to use BP to compute exact marginals on chains and trees. In this section, we discuss an algorithm to compute  $p(xq|xv)$  for any kind of graph.

We will explain the algorithm by example. Consider the DGM in Figure 20.3(a). This model, from (Koller and Friedman 2009), is a hypothetical model relating various variables pertaining to a typical student. The corresponding joint has the following form:

$$P(C,D, I, G, S, L, J,H) \quad (20.38) =$$

$$P(C)P(D|C)P(I)P(G|I,D)P(S|I)P(L|G)P(J|L,S)P(H|G,J) \quad (20.39)$$

Note that the forms of the CPDs do not matter, since all our calculations will be symbolic. However, for illustration purposes, we will assume all variables are binary.

Before proceeding, we convert our model to undirected form. This is not required, but it makes for a more unified presentation, since the resulting method can then be applied to both DGMs and UGMs (and, as we will see in Section 20.3.1, to a variety of other problems that have nothing to do with graphical models). Since the computational complexity of inference in DGMs and UGMs is, generally speaking, the same, nothing is lost in this transformation from a computational point of view. To convert the DGM to a UGM, we simply define a potential or factor for every CPD, yielding

$$p(C,D, I, G, S, L, J,H) \quad (20.40)$$

$$= \psi_C(C)\psi_D(D,C)\psi_I(I)\psi_G(G, I,D)\psi_S(S, I)\psi_L(L,G)\psi_J(J, L,$$

$$S)\psi_H(H,G, J) \quad (20.41)$$

Since all the potentials are locally normalized, since they are CPDs, there is no need for a global normalization constant, so  $Z = 1$ . The corresponding undirected graph is shown in Figure 20.3(b). Note that it has more edges than the DAG. In particular, any “unmarried” nodes that share a child must get “married”, by adding an edge between them; this process is known as moralization. Only then can the arrows be dropped. In this example, we added D-I, G-J, and L-S moralization arcs. The reason this operation is required is to ensure that the CI properties of the UGM match those of the DGM, as explained in Section 19.2.2. It also ensures there is a clique that can “store” the CPDs of each family.

Now suppose we want to compute  $p(J = 1)$ , the marginal probability that a person will get a job. Since we have 8 binary variables, we could simply enumerate over all possible assignments to all the variables (except for  $J$ ), adding up the probability of each joint instantiation:

$$p(J) = \sum_{LSGHID} p(C,D, I, G, S, L, J,H) \quad (20.42)$$

However, this would take  $O(2^8)$  time. We can be smarter by pushing sums inside products. This is the key idea behind the variable elimination algorithm (Zhang and Poole 1996), also called bucket elimination (Dechter 1996), or, in the context of genetic pedigree trees, the peeling algorithm (Cannings et al. 1978). In our example, we get

$$p(J) = \sum_{LSGHID} p(C,D, I, G, S, L, J,H)$$

$$= \psi_C(C)\psi_D(D,C)\psi_I(I)\psi_G(G, I,D)\psi_S(S, I)\psi_L(L,G) \\ \times \psi_J(J, L, S)\psi_H(H,G, J)$$

$$= LS \psi J(J, L, S) G \psi L(L, G) H \psi H(H, G, J) I \psi S(S, I) \psi I(I) \times D \psi G(G, I, D) C \psi C(C) \psi D(D, C)$$

We now evaluate this expression, working right to left as shown in Table 20.1. First we multiply together all the terms in the scope of the C operator to create the temporary factor

$$\tau 1(C, D) = \psi C(C) \psi D(D, C) \quad (20.43)$$

Then we marginalize out C to get the new factor

$$\tau 1(D) = C \tau 1(C, D) \quad (20.44)$$

Next we multiply together all the terms in the scope of the D operator and then marginalize out to create

$$\tau 2(G, I, D) = \psi G(G, I, D) \tau 1(D) \quad (20.45)$$

$$\tau 2(G, I) = D \tau 2(G, I, D) \quad (20.46)$$

Next we multiply together all the terms in the scope of the I operator and then marginalize out to create

$$\tau 3(G, I, S) = \psi S(S, I) \psi I(I) \tau 2(G, I) \quad (20.47)$$

$$\tau 3(G, S) = I \tau 3(G, I, S) \quad (20.48)$$

And so on. The above technique can be used to compute any marginal of interest, such as  $p(J)$  or  $p(J, H)$ . To compute a conditional, we can take a ratio of two marginals, where the visible variables have been clamped to their known values (and hence don't need to be summed over). For example,

$$p(J = j | I = 1, H = 0) = p(J = j, I = 1, H = 0) / p(J = j, I = 1, H = 0) \quad (20.49)$$

In general, we can write

$$p(xq|xv) = p(xq, xv) / p(xv) = xh p(xh, xq, xv) / xh x q p(xh, xq, xv) \quad (20.50)$$

The normalization constant in the denominator,  $p(xv)$ , is called the probability of the evidence.

See variableElimination for a simple Matlab implementation of this algorithm, which works for arbitrary graphs, and arbitrary discrete factors. But before you go too crazy, please read Section 20.3.2, which points out that VE can be exponentially slow in the worst case.

### 20.3.1 The generalized distributive law

Abstractly, VE can be thought of as computing the following expression:

$$p(xq|xv) \propto x c \psi c(xc) \quad (20.51)$$

It is understood that the visible variables  $xv$  are clamped, and not summed over. VE uses non-serial dynamic programming (Bertele and Brioschi 1972), caching intermediate results to avoid redundant computation.

However, there are other tasks we might like to solve for any given graphical model. For example, we might want the MAP estimate:

$$x^* = \text{argmax } x c \psi c(xc) \quad (20.52)$$

Fortunately, essentially the same algorithm can also be used to solve this task: we just replace sum with max. (We also need a traceback step, which actually recovers the argmax, as opposed to just the value of max; these details are explained in Section 17.4.4.)

In general, VE can be applied to any commutative semiring. This is a set  $K$ , together with two binary operations called “ $+$ ” and “ $\times$ ”, which satisfy the following three axioms:

1. The operation “ $+$ ” is associative and commutative, and there is an additive identity element called “0” such that  $k + 0 = k$  for all  $k \in K$ .
2. The operation “ $\times$ ” is associative and commutative, and there is a multiplicative identity element called “1” such that  $k \times 1 = k$  for all  $k \in K$ .
3. The distributive law holds, i.e.,

$$(a \times b) + (a \times c) = a \times (b + c) \quad (20.53)$$

for all triples  $(a, b, c)$  from  $K$ .

This framework covers an extremely wide range of important applications, including constraint satisfaction problems (Bistarelli et al. 1997; Dechter 2003), the fast Fourier transform (Aji and McEliece 2000), etc. See Table 20.2 for some examples.

### 20.3.2 Computational complexity of VE

The running time of VE is clearly exponential in the size of the largest factor, since we have sum over all of the corresponding variables. Some of the factors come from the original model (and are thus unavoidable), but new factors are created in the process of summing out. For example, in Equation 20.47, we created a factor involving  $G, I$  and  $S$ ; but these nodes were not originally present together in any factor. The order in which we perform the summation is known as the elimination order. This can have a large impact on the size of the intermediate factors that are created. For example, consider the ordering in Table 20.1: the largest created factor (beyond the original ones in the model) has size 3, corresponding to  $\tau_5(J, L, S)$ . Now consider the ordering in Table 20.3: now the largest factors are  $\tau_1(I, D, L, J, H)$  and  $\tau_2(D, L, S, J, H)$ , which are much bigger. We can determine the size of the largest factor graphically, without worrying about the actual numerical values of the factors. When we eliminate a variable  $X_t$ , we connect it to all variables that share a factor with  $X_t$  (to reflect the new temporary factor  $\tau_t$ ). The edges created by this process are called fill-in edges. For example, Figure 20.4 shows the fill-in edges introduced when we eliminate in the order  $C, D, I, \dots$ . The

first two steps do not introduce any fill-ins, but when we eliminate I, we connect G and S, since they co-occur in Equation 20.48. Let  $G(<)$  be the (undirected) graph induced by applying variable elimination to  $G$  using elimination ordering  $<$ . The temporary factors generated by VE correspond to maximal cliques in the graph  $G(<)$ . For example, with ordering (C,D, I,H, G, S, L), the maximal cliques are as follows:

$$\{C, D\}, \{D, I, G\}, \{G, L, S, J\}, \{G, J, H\}, \{G, I, S\} \quad (20.54)$$

It is clear that the time complexity of VE is

$$c \in C(G(<)) K^{|c|} \quad (20.55)$$

where  $C$  are the cliques that are created,  $|c|$  is the size of the clique  $c$ , and we assume for notational simplicity that all the variables have  $K$  states each. Let us define the induced width of a graph given elimination ordering  $<$ , denoted  $w(<)$ , as the size of the largest factor (i.e., the largest clique in the induced graph) minus 1. Then it is easy to see that the complexity of VE with ordering  $<$  is  $O(Kw(<)+1)$ . Obviously we would like to minimize the running time, and hence the induced width. Let us define the treewidth of a graph as the minimal induced width.

$$w = \min < \max c \in G(<) |c| - 1 \quad (20.56)$$

Then clearly the best possible running time for VE is  $O(DKw+1)$ . Unfortunately, one can show that for arbitrary graphs, finding an elimination ordering  $<$  that minimizes  $w(<)$  is NP-hard (Arnborg et al. 1987). In practice greedy search techniques are used to find reasonable orderings (Kjaerulff 1990), although people have tried other heuristic methods for discrete optimization, such as genetic algorithms (Larranaga et al. 1997). It is also possible to derive approximate algorithms with provable performance guarantees (Amir 2010). In some cases, the optimal elimination ordering is clear. For example, for chains, we should work forwards or backwards in time. For trees, we should work from the leaves to the root. These orderings do not introduce any fill-in edges, so  $w=1$ . Consequently, inference in chains and trees takes  $O(VK^2)$  time. This is one reason why Markov chains and Markov trees are so widely used. Unfortunately, for other graphs, the treewidth is large. For example, for an  $m \times n$  2d lattice, the treewidth is  $O(\min\{m,n\})$  (Lipton and Tarjan 1979). So VE on a  $100 \times 100$  Ising model would take  $O(2100)$  time. Of course, just because VE is slow doesn't mean that there isn't some smarter algorithm out there. We discuss this issue in Section 20.5.

### 20.3.3 A weakness of VE

The main disadvantage of the variable elimination

algorithm (apart from its exponential dependence on treewidth) is that it is inefficient if we want to compute multiple queries conditioned on the same evidence. For example, consider computing all the marginals in a chain-structured graphical model such as an HMM. We can easily compute the final marginal  $p(xT|v)$  by eliminating all the nodes  $x_1$  to  $x_{T-1}$  in order. This is equivalent to the forwards algorithm, and takes  $O(K^2T)$  time. But now suppose we want to compute  $p(x_{T-1}|v)$ . We have to run VE again, at a cost of  $O(K^2T)$  time. So the total cost to compute all the marginals is  $O(K^2T^2)$ . However, we know that we can solve this problem in  $O(K^2T)$  using forwards-backwards. The difference is that FB caches the messages computed on the forwards pass, so it can reuse them later. The same argument holds for BP on trees. For example, consider the 4-node tree in Figure 20.5. We can compute  $p(x_1|v)$  by eliminating  $x_2, x_3, x_4$ ; this is equivalent to sending messages up to  $x_1$  (the messages correspond to the  $\tau$  factors created by VE). Similarly we can compute  $p(x_2|v)$ ,  $p(x_3|v)$  and then  $p(x_4|v)$ . We see that some of the messages used to compute the marginal on one node can be re-used to compute the marginals on the other nodes. By storing the messages for later re-use, we can compute all the marginals in  $O(DK^2)$  time. This is what the up-down (collect-distribute) algorithm on trees does. The question is: how can we combine the efficiency of BP on trees with the generality of VE? The answer is given in Section 20.4.

## 翻译

**图20.3 左：“学生” DGM。右：等效的UGM。我们添加了道德化弧 D-I、G-J 和 L-S。**

**图20.4 消元过程的示例，顺序为 C、D、I 等。当我们消去 I 时（图 c），我们在 G 和 S 之间添加了一条填充边，因为它们没有直接连接。**

**图20.5 在树上发送多个消息。**

- (a)  $X_1$  是根节点。
- (b)  $X_2$  是根节点。
- (c)  $X_4$  是根节点。
- (d) 计算所有单一边际所需的所有消息。

**表20.1 按照顺序 C、D、I、H、G、S、L 从图20.3中消去变量，以计算 P(J)。**

**表20.2 一些交换半环。**

**表20.3 按照顺序 G、I、S、L、H、C、D 从图20.3中消去变量。**

## 20.3 变量消元算法

我们已经了解了如何使用BP算法在链和树上计算精确的边际概率。在本节中，我们将讨论一个算法，用于计算任意图中的  $p(x_q | x_v)$ 。

我们通过示例来解释这个算法。考虑图20.3(a)中的DGM模型。这个模型（来自Koller和Friedman 2009）是一个假设的模型，描述了一个典型学生的各种变量之间的关系。对应的联合分布具有以下形式：

$$P(C, D, I, G, S, L, J, H) = P(C)P(D|C)P(I)P(G|I, D)P(S|I)P(L|G)P(J|L, S)P(H|G, J)$$

请注意，CPD（条件概率分布）的形式不重要，因为所有计算都将是符号化的。然而，为了说明问题，我们假设所有变量都是二进制的。

在继续之前，我们将模型转换为无向图形式（UGM）。这不是必需的，但它使得演示更为统一，因为这样得到的方法可以应用于DGM和UGM（正如我们将在20.3.1节中看到的，它还可以应用于许多与图模型无关的其他问题）。由于DGM和UGM中的推断计算复杂度通常相同，从计算角度来看，这个转换不会丢失任何东西。

要将DGM转换为UGM，我们简单地为每个CPD定义一个潜力或因子，得到：

$$p(C, D, I, G, S, L, J, H) = \psi_C(C)\psi_D(D, C)\psi_I(I)\psi_G(G, I, D)\psi_S(S, I)\psi_L(L, G)\psi_J(J, L, S)\psi_H(H, G, J)$$

由于所有潜力都是局部归一化的（因为它们是CPD），因此不需要全局归一化常数，所以  $Z = 1$ 。对应的无向图如图20.3(b)所示。请注意，它比DAG（有向无环图）有更多的边。特别是，任何共享子节点的“未婚”节点必须通过添加边来“结婚”，这个过程称为道德化（moralization）。只有这样才能删除箭头。在这个例子中，我们添加了D-I、G-J和L-S的道德化弧。这一步操作的原因是确保UGM的条件独立性（CI）性质与DGM相匹配，正如在19.2.2节中解释的那样。它还确保存在一个可以“存储”每个家庭CPD的团。

现在假设我们想要计算  $p(J = 1)$ ，即一个人获得工作的边际概率。由于我们有8个二进制变量，我们可以简单地列举所有变量（除了J）可能的赋值，计算每个联合实例的概率：

$$p(J) = \sum_{L,S,G,H,I,D,C} p(C, D, I, G, S, L, J, H)$$

然而，这将需要  $O(2^7)$  时间。我们可以通过将求和推入乘积中变得更智能。这就是变量消元算法的关键思想（Zhang 和 Poole 1996），也称为桶消元（bucket elimination）（Dechter 1996），或者在遗传谱系树的背景下，叫做剥皮算法（peeling algorithm）（Cannings 等人 1978）。在我们的例子中，我们得到：

$$p(J) = \sum_{L,S,G,H,I,D,C} \psi_C(C)\psi_D(D, C)\psi_I(I)\psi_G(G, I, D)\psi_S(S, I)\psi_L(L, G)\psi_J(J, L, S)\psi_H(H, G, J)$$

我们现在从右到左计算这个表达式，如表20.1所示。首先，我们将C算子的所有项相乘，得到临时因子：

$$\tau_1(C, D) = \psi_C(C)\psi_D(D, C)$$

然后我们对C进行边际化，得到新的因子：

$$\tau_1(D) = \sum_C \tau_1(C, D)$$

接下来，我们将D算子的所有项相乘，然后进行边际化，得到：

$$\tau_2(G, I, D) = \psi_G(G, I, D)\tau_1(D)$$

然后对D进行边际化得到：

$$\tau_2(G, I) = \sum_D \tau_2(G, I, D)$$

接下来的步骤类似。上述技术可以用来计算任何感兴趣的边际，比如  $p(J)$  或  $p(J, H)$ 。

为了计算条件概率，我们可以取两个边际的比值，其中可见变量被固定为已知值（因此不需求和）。例如：

$$p(J = j | I = 1, H = 0) = \frac{p(J = j, I = 1, H = 0)}{\sum_j p(J = j, I = 1, H = 0)}$$

一般来说，我们可以写成：

$$p(x_q | x_v) = \frac{p(x_q, x_v)}{p(x_v)} = \sum_{x_h} \frac{p(x_h, x_q, x_v)}{p(x_h, x_q, x_v)}$$

分母中的归一化常数  $p(x_v)$  称为证据的概率。

## 解释

### 1. 变量消元 (Variable Elimination, VE) 算法：

- 这个算法用于计算图模型中的边际概率，尤其是对于一般图（不仅仅是树或链）很有效。通过将求和操作推入乘积中，VE算法大大减少了计算量，使得可以在不穷举所有可能的情况下计算概率。

### 2. 转化为无向图 (UGM)：

- 我们将有向图模型 (DGM) 转化为无向图模型 (UGM)，这使得算法更加统一，能同时适用于DGM和UGM。转化过程中使用了“道德化” (moralization) 技术：如果两个没有直接连接的节点共享一个子节点，它们之间就需要添加边，这样才能正确地表示条件独立性。

### 3. 消元过程：

- 变量消元通过逐个消去变量来简化计算，每次消除一个变量都会计算一个临时因子，并对该变量进行边际化。通过这种方式，逐步减少图中的变量，从而得到目标概率。

#### 4. 条件概率和边际概率：

- 消元算法不仅可以计算边际概率，还可以通过计算两个边际的比值来计算条件概率。

#### 5. 计算复杂度：

- 变量消元的时间复杂度是指数级的，因为每次消除变量时都会创建一个临时因子，计算中间结果的大小受变量消除顺序的影响。找到最优的消元顺序是一个NP难题，但可以使用贪心算法等方法来找到合理的顺序。

#### 6. VE的弱点：

- 变量消元算法在多次计算条件概率或边际概率时效率较低，因为每次查询都需要重新计算。这和前向后向算法或信念传播算法不同，它们通过缓存之前的计算结果来避免重复工作。

### 总结要点：

- 变量消元（VE）算法是一种用于计算图模型中边际概率和条件概率的算法，通过消去变量并逐步简化计算过程。
- 将有向图转化为无向图，并通过道德化确保条件独立性匹配。
- VE的计算过程中，通过逐个消去变量并计算临时因子来逐步得到目标概率。
- 变量消元算法的计算复杂度是指数级的，且依赖于消元顺序。为了提高效率，可以使用贪心策略等方法找到合理的顺序。
- VE在多次查询时效率较低，这个问题在前向后向算法或信念传播算法中通过缓存结果得到了更好的解决。

这一节主要在讲什么？我不太懂

这一节主要讲的是**变量消元算法 (Variable Elimination, VE)**，它是一种用于在图模型（比如概率图模型）中计算边际概率和条件概率的算法。

### 主要内容概述：

#### 1. 变量消元算法 (VE)：

- 变量消元是一种用于计算图模型中边际概率和条件概率的技术，特别是当图是一个复杂的结构时。它通过逐步消去图中的变量，并使用某些数学操作（如求和和乘积）来简化计算，避免穷举所有可能的变量组合。
- 简单来说，就是通过逐步“消除”不重要的变量，最终计算出我们感兴趣的概率。

#### 2. 过程示例：

- 在本节中，作者给出了一个示例（假设的“学生模型”），通过这个模型来展示如何用变量消元算法计算某些概率（比如某人找到工作的概率）。这部分解释了如何通过逐步消去变量并结合相应的因子（如条件概率分布）来得到目标概率。

### 3. 道德化 (Moralization) :

- 在一些情况下，作者将有向图转换为无向图 (UGM)，以便更容易应用变量消元算法。这一过程中涉及了“道德化”的概念，即如果两个没有直接连接的节点有共同的子节点，它们之间必须添加边。这是为了确保图中每一对相互独立的节点都符合模型的条件独立性。

### 4. 边际和条件概率：

- 变量消元不仅可以用来计算**边际概率**（比如“这个人找到工作的概率”），还可以通过计算两个边际的比值来计算**条件概率**（比如“在某些条件下，这个人找到工作的概率”）。

### 5. 计算复杂度：

- 变量消元算法的效率是一个重要问题，特别是在处理复杂图时。算法的运行时间是指数级的，主要受**消元顺序**（消去哪些变量的顺序）影响。选择好的消元顺序可以大大减少中间计算的复杂度。

### 6. VE的弱点：

- 如果你需要对同一证据进行多次查询，变量消元会很高效，因为每次查询都会重新计算一次，这不同于前向后向算法或信念传播算法，它们会存储计算结果并加速后续的计算。

## 总结：

这一节讲解了变量消元算法，它通过逐步消去变量来计算边际和条件概率。算法涉及一些技术细节，如道德化过程和消元顺序的选择。尽管VE算法在处理某些图模型时非常有效，但它的计算复杂度很高，尤其是在多次查询的情况下。因此，在某些情况下，其他算法（如前向后向或信念传播）可能更加高效。

希望这个解释能帮助你理解这一节的内容！如果你有任何问题或者不明白的地方，请随时问我！

翻译并解释以下内容，翻译不能有遗漏，解释请通俗易懂且详尽，因为我的数学基础不太好，最后总结要点。

Figure 20.6 Left: this graph is not triangulated, despite appearances, since it contains a chordless 5-cycle 1-2-3-4-5-1. Right: one possible triangulation, by adding the 1-3 and 1-4 fill-in edges.

Figure 20.7 (a) The student graph with fill-in edges added.

(b) The maximal cliques. (c) The junction tree. An edge between nodes s and t is labeled by the intersection of the sets on nodes s and t; this is called the separating set. Figure 20.8 The junction tree derived from an HMM/SSM of length T =4.

## 20.4 The junction tree algorithm

The junction tree algorithm or JTA generalizes BP from trees to arbitrary graphs. We sketch the basic idea below; for details, see e.g., (Koller and Friedman 2009).

### 20.4.1 Creating a junction tree

The basic idea behind the JTA is this. We first run the VE algorithm “symbolically”, adding fill-in edges as we go, according to a given elimination ordering. The resulting graph will be a chordal graph, which means that every undirected cycle  $X_1 - X_2 \dots X_k - X_1$  of length  $k \geq 4$  has a chord, i.e., an edge connects  $X_i, X_j$  for all non-adjacent nodes  $i, j$  in the cycle.<sup>2</sup> Having created a chordal graph, we can extract its maximal cliques. In general, finding max cliques is computationally hard, but it turns out that it can be done efficiently from this special kind of graph. Figure 20.7(b) gives an example, where the max cliques are as follows: {C,D}, {G, I,D}, {G, S, I}, {G,J, S,L}, {H,G, J} (20.57)

Note that if the original graphical model was already chordal, the elimination process would not add any extra fill-in edges (assuming the optimal elimination ordering was used). We call such models decomposable, since they break into little pieces defined by the cliques. It turns out that the cliques of a chordal graph can be arranged into a special kind of tree known as a junction tree. This enjoys the running intersection property (RIP), which means that any subset of nodes containing a given variable forms a connected component. Figure 20.7(c) gives an example of such a tree. We see that the node I occurs in two adjacent tree nodes, so they can share information about this variable. A similar situation holds for all the other variables. One can show that if a tree that satisfies the running intersection property, then applying BP to this tree (as we explain below) will return the exact values of  $p(x_c|v)$  for each node c in the tree (i.e., clique in the induced graph). From this, we can easily extract the node and edge marginals,  $p(x_t|v)$  and  $p(x_s, x_t|v)$  from the original model, by marginalizing the clique distributions.

### 20.4.2 Message passing on a junction tree

Having constructed a junction tree, we can use it for inference. The process is very similar to belief propagation on a tree. As in Section 20.2, there are two versions: the

sum-product form, also known as the Shafer-Shenoy algorithm, named after (Shafer and Shenoy 1990); and the belief updating form (which involves division), also known as the Hugin (named after a company) or the Lauritzen-Spiegelhalter algorithm (named after (Lauritzen and Spiegelhalter 1988)). See (Lepar and Shenoy 1998) for a detailed comparison of these methods. Below we sketch how the Hugin algorithm works. We assume the original model has the following form:

$$p(x) = \frac{1}{Z} \prod_{c \in C(G)} \psi_c(x_c) \quad (20.58)$$

where  $C(G)$  are the cliques of the original graph. On the other hand, the tree defines a distribution of the following form:

$$p(x) = \prod_{c \in C(T)} \psi_c(x_c) \prod_{s \in S(T)} \psi_s(x_s) \quad (20.59)$$

where  $C(T)$  are the nodes of the junction tree (which are the cliques of the chordal graph), and  $S(T)$  are the separators of the tree. To make these equal, we initialize by defining  $\psi_s = 1$  for all separators and  $\psi_c = 1$  for all cliques. Then, for each clique in the original model,  $c \in C(G)$ , we find a clique in the tree  $c' \in C(T)$  which contains it,  $c \supseteq c'$ . We then multiply  $\psi_c$  onto  $\psi_{c'}$  by computing  $\psi_{c'} = \psi_c \psi_{c'}$ . After doing this for all the cliques in the original graph, we have

$$\prod_{c \in C(T)} \psi_c(x_c) = \prod_{c \in C(G)} \psi_c(x_c) \quad (20.60)$$

As in Section 20.2.1, we now send messages from the leaves to the root and back, as sketched in Figure 20.1. In the upwards pass, also known as the collect-to-root phase, node  $i$  sends to its parent  $j$  the following message:

$$m_{i \rightarrow j}(S_{ij}) = \psi_j(S_{ij}) \psi_i(C_i) \quad (20.61)$$

That is, we marginalize out the variables that node  $i$  "knows about" which are irrelevant to  $j$ , and then we send what is left over. Once a node has received messages from all its children, it updates its belief state using

$$\psi_i(C_i) \propto \psi_i(C_i) \prod_{j \in \text{children } i} m_{j \rightarrow i}(S_{ij}) \quad (20.62)$$

At the root,  $\psi_r(C_r)$  represents  $p(x_r | v)$ , which is the posterior over the nodes in clique  $C_r$  conditioned on all the evidence. Its normalization constant is  $p(v)/Z_0$ , where  $Z_0$  is the normalization constant for the unconditional prior,  $p(x)$ . (We have  $Z_0 = 1$  if the original model was a DGM.) In the downwards pass, also known as the distribute-from-root phase, node  $i$  sends to its children  $j$  the following message:

$$m_{i \rightarrow j}(S_{ij}) = \psi_j(S_{ij}) \psi_i(C_i) m_{r \rightarrow i}(S_{ri}) \quad (20.63)$$

We divide out by what  $j$  sent to  $i$  to avoid double counting the evidence. This requires that we store the messages from the upwards pass. Once a node has received a top-down message from its parent, it can compute its final

belief state using

$$\psi_j(C_j) \propto \psi_j(C_j) m_i \rightarrow j(S_{ij}) \quad (20.64)$$

An equivalent way to present this algorithm is based on storing the messages inside the separator potentials. So on the way up, sending from  $i$  to  $j$  we compute the separator potential

$$\psi^*_{ij}(S_{ij}) = C_i \setminus S_{ij} \psi_i(C_i) \quad (20.65)$$

and then update the recipient potential:

$$\psi^*_{ij}(C_j) \propto \psi_j(C_j) \psi^*_{ij}(S_{ij}) \psi_{ij}(S_{ij}) \quad (20.66)$$

(Recall that we initialize  $\psi_{ij}(S_{ij}) = 1$ .) This is sometimes called passing a flow from  $i$  to  $j$ . On the way down, from  $i$  to  $j$ , we compute the separator potential

$$\psi^{**}_{ij}(S_{ij}) = C_i \setminus S_{ij} \psi^*_{ij}(C_i) \quad (20.67)$$

and then update the recipient potential:

$$\psi^{**}_{ij}(C_j) \propto \psi^*_{ij}(C_j) \psi^{**}_{ij}(S_{ij}) \psi^*_{ij}(S_{ij}) \quad (20.68)$$

This process is known as junction tree calibration. See Figure 20.1 for an illustration. Its correctness follows from the fact that each edge partitions the evidence into two distinct groups, plus the fact that the tree satisfies RIP, which ensures that no information is lost by only performing local computations.

#### 20.4.2.1 Example: jtreet algorithm on a chain

It is interesting to see what happens if we apply this process to a chain structured graph such as an HMM. A detailed discussion can be found in (Smyth et al. 1997), but the basic idea is this. The cliques are the edges, and the separators are the nodes, as shown in Figure 20.8. We initialize the potentials as follows: we set  $\psi_s = 1$  for all the separators, we set  $\psi_c(xt-1, xt) = p(xt|xt-1)$  for clique  $c = (Xt-1, Xt)$ , and we set  $\psi_c(xt, yt) = p(yt|xt)$  for clique  $c = (Xt, Yt)$ .

Next we send messages from left to right. Consider clique  $(Xt-1, Xt)$  with potential  $p(Xt|Xt-1)$ . It receives a message from clique  $(Xt-2, Xt-1)$  via separator  $Xt-1$  of the form  $x_{t-2} p(Xt-2, Xt-1 | v_{1:t-1}) = p(Xt-1 | v_{1:t-1})$ . When combined with the clique potential, this becomes the two-slice predictive density

$$p(Xt|Xt-1)p(Xt-1 | v_{1:t-1}) = p(Xt-1, Xt | v_{1:t-1}) \quad (20.69)$$

The clique  $(Xt-1, Xt)$  also receives a message from  $(Xt, Yt)$  via separator  $Xt$  of the form  $p(yt|Xt)$ , which corresponds to its local evidence. When combined with the updated clique potential, this becomes the two-slice filtered posterior

$$p(Xt-1, Xt | v_{1:t-1})p(yt|Xt) = p(Xt-1, Xt | v_{1:t}) \quad (20.70)$$

Thus the messages in the forwards pass are the filtered belief states  $x_t$ , and the clique potentials are the two-slice distributions. In the backwards pass, the messages are the

update factors  $\gamma_t$  at , where  $\gamma_t(k) = p(x_t = k | v_1:T)$  and  $\alpha_t(k) = p(x_t = k | v_1:t)$ . By multiplying by this message, we “swap out” the old  $\alpha_t$  message and “swap in” the new  $\gamma_t$  message. We see that the backwards pass involves working with posterior beliefs, not conditional likelihoods. See Section 18.3.2.3 for further discussion of this difference.

#### 20.4.3 Computational complexity of JTA

If all nodes are discrete with  $K$  states each, it is clear that the JTA takes  $O(|C|Kw+1)$  time and space, where  $|C|$  is the number of cliques and  $w$  is the treewidth of the graph, i.e., the size of the largest clique minus 1. Unfortunately, choosing a triangulation so as to minimize the treewidth is NP-hard, as explained in Section 20.3.2. The JTA can be modified to handle the case of Gaussian graphical models. The graph-theoretic steps remain unchanged. Only the message computation differs. We just need to define how to multiply, divide, and marginalize Gaussian potential functions. This is most easily done in information form. See e.g., (Lauritzen 1992; Murphy 1998; Cemgil 2001) for the details. The algorithm takes  $O(|C|w^3)$  time and  $O(|C|w^2)$  space. When applied to a chain structured graph, the algorithm is equivalent to the Kalman smoother in Section 18.3.2.

#### 20.4.4 JTA generalizations

We have seen how to use the JTA algorithm to compute posterior marginals in a graphical model. There are several possible generalizations of this algorithm, some of which we mention below. All of these exploit graph decomposition in some form or other. They only differ in terms of how they define/ compute messages and “beliefs”. The key requirement is that the operators which compute messages form a commutative semiring (see Section 20.3.1).

- Computing the MAP estimate. We just replace the sum-product with max-product in the collect phase, and use traceback in the distribute phase, as in the Viterbi algorithm (Section 17.4.4). See (Dawid 1992) for details.
- Computing the N-most probable configurations (Nilsson 1998).
- Computing posterior samples. The collect pass is the same as usual, but in the distribute pass, we sample variables given the values higher up in the tree, thus generalizing forwardsfiltering backwards-sampling for HMMs described in Section 17.4.5. See (Dawid 1992) for details.
- Solving constraint satisfaction problems (Dechter 2003).

- Solving logical reasoning problems (Amir and McIlraith 2005).

## 翻译：

**图20.6 左：**这个图尽管看起来是三角化的，但实际上并不是，因为它包含一个没有弦的5-cycle (1-2-3-4-5-1)。右：一种可能的三角化方法，通过添加1-3和1-4的填充边。

**图20.7 (a)** 学生图，已添加填充边。**(b)** 最大团。**(c)** 联接树。节点s和t之间的边被标记为s和t节点集合的交集，这称为分隔集。

**图20.8** 从T=4的HMM/SSM中导出的联接树。

## 20.4 联接树算法 (JTA)

### 20.4.1 创建联接树

联接树算法 (JTA) 是从树状图到任意图的推断方法的推广。下面简要描述其基本思想；详细内容请参考如 (Koller 和 Friedman 2009)。

JTA的基本思路是：首先，我们按给定的消元顺序，使用变量消元算法 (VE) “符号化”地进行计算，同时添加填充边。最终得到的图会是一个**和弦图**，这意味着图中每个长度大于等于4的无向环，必须有一条弦连接每对非邻接节点 i 和 j。创建好和弦图后，我们可以提取其中的最大团。一般来说，找到最大团是一个计算上困难的问题，但对于这种特殊类型的图，计算可以高效进行。图20.7(b)给出了一个例子，其中最大团如下： $\{C, D\}$ ,  $\{G, I, D\}$ ,  $\{G, S, I\}$ ,  $\{G, J, S, L\}$ ,  $\{H, G, J\}$  (式20.57)。

如果原始图已经是和弦图，那么消元过程不会添加额外的填充边（假设使用了最佳消元顺序）。我们称这种模型为**可分解模型**，因为它们可以被分解成由团定义的小块。和弦图的团可以组成一种特殊类型的树，称为**联接树**。这种树具有**运行交集属性 (RIP)**，这意味着包含给定变量的任何节点子集都构成一个连通组件。图20.7(c)给出了这样一棵树的示例。我们看到节点I出现在两个相邻的树节点中，因此它们可以共享该变量的信息。类似地，其他变量也具有类似的共享情况。如果一棵树满足运行交集属性，那么在这棵树上应用信念传播 (BP)，将返回树中每个节点的精确后验值（即所生成图中的团）。从这里，我们可以通过边际化团分布轻松提取原始模型中的节点和边缘边际， $p(x_t|v)$  和  $p(x_s, x_t|v)$ 。

### 20.4.2 联接树上的消息传递

创建了联接树后，我们可以使用它来进行推断。这个过程非常类似于树上的信念传播。在第20.2节中，我们介绍了两种版本：和积形式，也称为Shafer-Shenoy算法（命名自Shafer和Shenoy 1990）；和信念更新形式（涉及除法），也称为Hugin算法（命名自某公司）或Lauritzen-Spiegelhalter算法（命名自Lauritzen和Spiegelhalter 1988）。有关这些方法的详细比较，请参见 (Lepar和Shenoy 1998)。以下是Hugin算法的工作原理。

我们假设原始模型具有以下形式：

$$p(x) = \frac{1}{Z} \prod_{c \in C(G)} \psi_c(x_c)$$

其中， $C(G)$  是原始图的团。

另一方面，树定义了以下形式的分布：

$$p(x) = \prod_{c \in C(T)} \psi_c(x_c) \prod_{s \in S(T)} \psi_s(x_s)$$

其中， $C(T)$  是联接树的节点（即和弦图的团）， $S(T)$  是树的分隔符。为了使这两者相等，我们通过初始化  $\psi_s = 1$  对所有分隔符，和初始化  $\psi_c = 1$  对所有团。然后，对于原始模型中的每个团  $c \in C(G)$ ，我们找到联接树中的一个包含它的团  $c \in C(T)$ ，即  $c \subseteq c$ 。接着我们将  $\psi_c$  乘到  $\psi_c$  上，计算  $\psi_c = \psi_c \times \psi_c$ 。完成所有原始图的团后，我们有：

$$\prod_{c \in C(T)} \psi_c(x_c) = \prod_{c \in C(G)} \psi_c(x_c)$$

接下来，按照第20.2.1节的方式，我们从叶节点向根节点发送消息，并返回。上行过程，也称为**收集到根节点阶段**，节点  $i$  向父节点  $j$  发送以下消息：

$$m_{i \rightarrow j}(S_{ij}) = \sum_{C_i \setminus S_{ij}} \psi_i(C_i)$$

即，我们对节点  $i$  已知的但与节点  $j$  无关的变量进行边际化，然后将剩下的部分发送出去。一旦节点收到来自所有子节点的消息，它就使用以下公式更新自己的信念状态：

$$\psi_i(C_i) \propto \psi_i(C_i) \prod_{j \in \text{ch}(i)} m_{j \rightarrow i}(S_{ij})$$

在根节点， $\psi_r(C_r)$  代表  $p(x_{C_r} | v)$ ，即在所有证据的条件下，根节点的后验概率。其归一化常数是  $p(v)/Z_0$ ，其中  $Z_0$  是无条件先验的归一化常数（如果原始模型是DGM，则  $Z_0 = 1$ ）。

在下行过程，也称为**从根节点分发阶段**，节点  $i$  向其子节点  $j$  发送以下消息：

$$m_{i \rightarrow j}(S_{ij}) = \frac{\sum_{C_i \setminus S_{ij}} \psi_i(C_i)}{m_{j \rightarrow i}(S_{ij})}$$

我们通过除去节点  $j$  发送给  $i$  的消息，来避免重复计算证据。这个过程需要存储来自上行阶段的消息。收到来自父节点的上行消息后，节点可以使用以下公式计算最终的信念状态：

$$\psi_j(C_j) \propto \psi_j(C_j) \times m_{i \rightarrow j}(S_{ij})$$

这种算法的另一种表现方式是将消息存储在分隔符潜力中。上行过程中，计算分隔符潜力：

$$\psi_{ij}^*(S_{ij}) = \sum_{C_i \setminus S_{ij}} \psi_i(C_i)$$

然后更新接收节点的潜力：

$$\psi_j^*(C_j) \propto \psi_j(C_j) \times \psi_{ij}^*(S_{ij}) \times \psi_{ij}(S_{ij})$$

下行过程中，计算更新后的分隔符潜力：

$$\psi_{ij}^{**}(S_{ij}) = \sum_{C_i \setminus S_{ij}} \psi_i^*(C_i)$$

然后更新接收节点的潜力：

$$\psi_j^{**}(C_j) \propto \psi_j^*(C_j) \times \psi_{ij}^{**}(S_{ij}) \times \psi_{ij}^*(S_{ij})$$

这个过程称为**联接树校准**。

## 总结要点：

- 联接树算法 (JTA)**：JTA是一种从树到任意图的推断算法。它通过将图转换为和弦图，并创建联接树，进而实现精确的推断。
- 创建联接树**：通过使用变量消元 (VE) 符号化地进行计算并添加填充边，我们得到一个和弦图。然后，从这个图中提取最大团，并用它们构建一个满足运行交集属性 (RIP) 的联接树。
- 消息传递**：一旦联接树构建完成，我们就可以像树上的信念传播一样，通过消息传递来进行推断。上行过程（收集到根节点）和下行过程（从根节点分发）涉及消息的更新和计算信念状态。
- 计算复杂度**：JTA的时间和空间复杂度是  $O(|C|K^{w+1})$ ，其中  $|C|$  是团的数量， $w$  是图的树宽度。尽管选择最佳三角化和最小化树宽度是NP难的，但实际应用中通常使用贪心方法来处理。
- 应用**：JTA不仅适用于图模型的边际计算，还可以进行最大后验 (MAP) 估计、后验采样等任务。

这一节主要在讲什么？我不太懂

这一节主要讲的是**联接树算法 (Junction Tree Algorithm, JTA)**，它是一种在图模型中进行推断的高级方法。特别地，JTA是对传统的信念传播 (BP) 算法的推广，可以应用于更复杂的图形结构，甚至是带有循环的图。

## 主要内容概述：

## 1. 联接树算法 (JTA) :

- JTA通过先将原始的图转换为一种叫做**和弦图**的特殊图，然后基于这个和弦图创建一个**联接树**，用来进行推断。它通过这种结构来处理图模型中的信息传播，解决原始图模型中的复杂计算问题。

## 2. 创建联接树的步骤：

- 首先，我们使用\*\*变量消元 (VE) \*\*算法来符号化地处理图模型，并通过给定的消元顺序来逐步消去变量，添加“填充边”。
- 通过这种消元和添加填充边的过程，我们得到一个和弦图。和弦图的特点是，图中的所有环都会有“弦”（即连接非相邻节点的边）。
- 接着，我们从和弦图中提取最大团 (cliques) ——这些最大团就是联接树的“节点”。
- 联接树中的节点（即最大团）之间通过\*\*分隔符 (separators) \*\*连接，这些分隔符帮助我们确保信息可以在树结构的不同部分之间流动。

## 3. 消息传递：

- 一旦创建了联接树，就可以像信念传播 (BP) 一样通过树上的节点之间传递消息来进行推断。
- 过程分为两个阶段：
  - **上行过程 (收集到根节点阶段)**：从树的叶节点开始，向根节点传递信息。
  - **下行过程 (从根节点分发阶段)**：从根节点向下传递信息，更新每个节点的信念状态。

## 4. 适用于更多问题：

- JTA不仅仅用于计算边际概率，还可以用于**最大后验 (MAP) 估计、后验采样**等任务。
- JTA的适用范围非常广泛，不仅限于标准的图模型，还可以用于其他一些复杂的推断任务。

## 5. 计算复杂度：

- JTA的计算复杂度依赖于图的结构，特别是图的**树宽度 (treewidth)**。在最坏情况下，选择最优的消元顺序是一个NP难的问题，但实际上可以通过一些启发式方法来找到合理的顺序，以避免计算量过大。

## 总结：

1. \*\*联接树算法 (JTA) \*\*是一种基于构建和弦图并生成联接树的推断算法，用于计算图模型中的边际概率、条件概率等。
2. 通过消元和填充边的过程，我们得到一个和弦图，并从中提取出最大团构建联接树。联接树用于高效地传递信息和进行推断。
3. JTA可以用来计算最大后验 (MAP) 估计、后验采样等任务，且适用于多种问题。
4. JTA的时间复杂度依赖于树宽度，最坏情况下计算复杂度很高，但可以通过启发式方法优化。

翻译并解释以下内容，翻译不能有遗漏，解释请通俗易懂且详尽，因为我的数学基础不太好，最后总结要点。

Figure 20.9 Encoding a 3-SAT problem on  $n$  variables and  $m$  clauses as a DGM. The  $Q_s$  variables are binary random variables. The  $C_t$  variables are deterministic functions of the  $Q_s$ 's, and compute the truth value of each clause. The  $A_t$  nodes are a chain of AND gates, to ensure that the CPT for the final  $x$  node has bounded size. The double rings denote nodes with deterministic CPDs.

Table 20.4 Summary of some methods that can be used for inference in graphical models. “D” means that all the hidden variables must be discrete. “L-G” means that all the factors must be linear-Gaussian. The term “single query” refers to the restriction that VE only computes one marginal  $p(x_q|x_v)$  at a time. See Section 20.3.3 for a discussion of this point. “C-E” stands for “conjugate exponential”; this means that variational mean field only applies to models where the likelihood is in the exponential family, and the prior is conjugate. This includes the D and LG case, but many others as well, as we will see in Section 21.5.

## 20.5 Computational intractability of exact inference in the worst case

As we saw in Sections 20.3.2 and 20.4.3, VE and JTA take time that is exponential in the treewidth of a graph. Since the treewidth can be  $O(\text{number of nodes})$  in the worst case, this means these algorithms can be exponential in the problem size.

Of course, just because VE and JTA are slow doesn't mean that there isn't some smarter algorithm out there. Unfortunately, this seems unlikely, since it is easy to show that exact inference is NP-hard (Dagum and Luby 1993). The proof is a simple reduction from the satisfiability problem. In particular, note that we can encode any 3-SAT problem<sup>4</sup> as a DGM with deterministic links, as shown in Figure 20.9. We clamp the final node,  $x$ , to be on, and we arrange the CPTs so that  $p(x=1) > 0$  iff there is a satisfying assignment. Computing any posterior marginal requires evaluating the normalization constant  $p(x=1)$ , which represents the probability of the evidence, so inference in this model implicitly solves the SAT problem. In fact, exact

inference is #P-hard (Roth 1996), which is even harder than NP-hard. (See e.g., (Arora and Barak 2009) for definitions of these terms.) The intuitive reason for this is that to compute the normalizing constant  $Z$ , we have to count how many satisfying assignments there are. By contrast, MAP estimation is provably easier for some model classes (Greig et al. 1989), since, intuitively speaking, it only requires finding one satisfying assignment, not counting all of them.

### 20.5.1 Approximate inference

Many popular probabilistic models support efficient exact inference, since they are based on chains, trees or low treewidth graphs. But there are many other models for which exact inference is intractable. In fact, even simple two node models of the form  $\theta \rightarrow x$  may not support exact inference if the prior on  $\theta$  is not conjugate to the likelihood  $p(x|\theta)$ .

Therefore we will need to turn to approximate inference methods. See Table 20.4 for a summary of coming attractions. For the most part, these methods do not come with any guarantee as to their accuracy or running time. Theoretical computer scientists would therefore describe them as heuristics rather than approximation algorithms. In fact, one can prove that it is not possible to construct polynomial time approximation schemes for inference in general discrete GMs (Dagum and Luby 1993; Roth 1996). Fortunately, we will see that for many of these heuristic methods often perform well in practice.

### 翻译：

**图20.9** 将一个含有 $n$ 个变量和 $m$ 个子句的3-SAT问题编码为DGM。 $Q_s$  变量是二元随机变量。 $C_t$  变量是  $Q_s$  的确定性函数，计算每个子句的真值。 $A_t$  节点是一个AND门链，用来确保最终  $x$  节点的条件概率表 (CPT) 大小是有限的。双环表示具有确定性CPD的节点。

**表20.4** 一些可以用于图模型推断的方法汇总。“D”表示所有的隐藏变量必须是离散的。“L-G”表示所有因子必须是线性高斯的。术语“单一查询”指的是变量消元 (VE) 仅计算一个边际概率  $p(x_q|x_v)$  的限制。见20.3.3节讨论此点。“C-E”代表“共轭指数” (conjugate exponential)；这意味着变分均值场仅适用于似然函数属于指数族且先验是共轭的模型。这包括D和L-G情况，但也包括许多其他情况，如我们将在21.5节中看到的。

## 20.5 最坏情况下精确推断的计算不可行性

正如我们在20.3.2节和20.4.3节中看到的，**变量消元 (VE) 和联接树算法 (JTA) 的计算时间是与图的树宽度成指数关系的。由于在最坏的情况下树宽度可以达到节点数的数量级**，这意味着这些算法的时间复杂度在最坏情况下可能会是问题规模的指级别。

然而，仅仅因为VE和JTA很慢，并不意味着不存在更智能的算法。遗憾的是，这似乎不太可能实现，因为很容易证明精确推断是NP难的（Dagum和Luby 1993）。这个证明是通过将可满足性问题（SAT问题）归约到精确推断来完成的。具体来说，注意我们可以将任何3-SAT问题编码为一个带有确定性连接的DGM，如图20.9所示。我们将最终节点 $x$ 固定为1，并安排条件概率表（CPT）使得 $p(x=1) > 0$ 当且仅当存在满足的赋值。计算任何后验边际都需要评估归一化常数 $p(x=1)$ ，它表示证据的概率，因此在此模型中的推断本质上解决了SAT问题。实际上，精确推断是\*\*#P难\*\*的（Roth 1996），这比NP难还要难。（参见例如（Arora和Barak 2009）以了解这些术语的定义。）这背后的直观原因是，为了计算归一化常数 $Z$ ，我们必须计数有多少个满足的赋值。相比之下，最大后验估计（MAP）对于某些模型类别来说被证明是更容易的（Greig等人 1989），因为直观地说，它只需要找到一个满足的赋值，而不是计算所有可能的赋值。

---

### 20.5.1 近似推断

许多流行的概率模型支持高效的精确推断，因为它们基于链、树或低树宽度的图形。但是，也有许多其他模型在精确推断上是不可行的。实际上，即使是简单的两个节点模型（例如 $\theta \rightarrow x$ ）如果 $\theta$ 的先验与似然 $p(x|\theta)$ 不是共轭的，也可能不支持精确推断。

因此，我们需要转向**近似推断方法**。请参见表20.4，了解未来要介绍的方法。大多数情况下，这些方法并没有提供关于准确性或运行时间的任何保证。因此，理论计算机科学家通常将它们描述为**启发式方法**，而不是近似算法。实际上，我们可以证明，对于一般的离散图模型，构建多项式时间的近似方案是不可能的（Dagum和Luby 1993；Roth 1996）。幸运的是，我们会看到，许多启发式方法在实践中通常能表现良好。

---

#### 解释：

这节主要讲的是**精确推断**和**近似推断**在图模型中的计算难度和方法。

##### 1. 精确推断的困难：

- 在图模型中，推断任务往往会变得非常复杂，尤其是在图的\*\*树宽度（treewidth）\*\*很大的情况下。**变量消元（VE）和联接树算法（JTA）的计算时间是指指数级的**，这意味着当问题规模（例如图的节点数）增大时，计算所需的时间会急剧增加。
- 更重要的是，推断问题本质上是NP难的，甚至比NP难度更高。通过将经典的3-SAT问题（一种著名的计算复杂问题）归约到图模型推断中，证明了精确推断的难度。

##### 2. 近似推断方法：

- 由于精确推断在很多情况下计算上不可行，研究者们转向了**近似推断方法**。近似推断不保证精确的结果，但在很多实际问题中表现得相当好。

- **近似推断方法**通常使用启发式算法，这些算法通常无法提供关于精确度和运行时间的理论保证，因此它们被认为是启发式方法，而不是标准的近似算法。

### 3. 启发式方法的实践效果：

- 尽管近似推断方法在理论上无法提供最优保证，但它们在很多实际问题中能够提供有效的解决方案。因此，虽然没有数学保证，但这些方法通常在实践中表现得很不错。

## 总结要点：

1. **精确推断的计算不可行性**：由于树宽度可能非常大，精确推断（例如通过VE和JTA）通常是不可行的，计算时间呈指数级增长。
2. **精确推断是NP难的**：通过将3-SAT问题归约到图模型推断，证明了推断是NP难的，甚至是更难的\*\*#P难\*\*。
3. **近似推断方法**：对于大多数模型，精确推断不可行，因此需要使用近似推断方法，这些方法通常表现良好，但无法提供理论上的准确性或时间保证。
4. **启发式方法**：许多近似推断方法是启发式的，虽然不能保证最优解，但它们通常能在实践中提供有效的解决方案。

这一节主要在讲什么？我不太懂

这一节主要讲的是图模型推断中的计算难度以及近似推断方法。

## 主要内容概述：

### 1. 推断的计算复杂性：

- 精确推断，如变量消元（VE）和联接树算法（JTA），在一些图模型中计算量是指数级增长的。这是因为这些算法的运行时间依赖于图的**树宽度**，在最坏情况下，树宽度可能和图的节点数一样大，从而导致计算非常慢。
- 精确推断是NP难的，也就是说，推断问题本质上是非常难解决的，计算量巨大。这部分通过将3-SAT问题（一种已知的计算难题）转化为图模型中的推断问题来证明。

### 2. 近似推断方法：

- 由于很多图模型的推断问题计算上不可行，我们需要使用**近似推断方法**。这些方法通常无法保证准确性或运行时间，但在实际应用中，它们表现得很好。
- 这些近似方法通常使用**启发式算法**，这种方法依靠经验或规则来找到近似解，而不是通过传统的精确算法。

### 3. 计算上不可行的原因：

- 精确推断之所以不可行，部分原因是我们需要计算归一化常数，这涉及到计算所有可能的满足条件的赋值，导致计算量非常大。相比之下，像\*\*最大后验估计 (MAP) \*\*这种方法只需要找到一个满足条件的赋值，计算上要容易得多。

## 总结要点：

- 精确推断的困难：**如变量消元和联接树算法在计算时可能呈指数级增长，尤其在图的树宽度较大的时候，推断问题非常难解决。
- NP难度：**通过将3-SAT问题归约到图模型推断，证明了推断是NP难的，甚至更复杂。
- 近似推断方法：**为了应对计算困难，使用近似推断方法，这些方法在实际应用中表现良好，但无法提供准确性和时间的理论保证。
- 启发式方法：**这些近似方法通常基于启发式算法，能在很多实际问题中提供有效解，但不能保证最优。

基于本对话里我贴过的整章原文（20.1-20.5），输出：①知识点总结（分点+关键词+用途）②子章节关联（依赖/对比/递进+关系句，附Mermaid可选）③本章总结。

## ① 知识点总结

- 信念传播 (Belief Propagation, BP)**
  - 关键词：**树结构图、消息传递、边际概率、后验概率、树上的推断
  - 用途：**BP主要用于链式结构或树结构图模型中，推断隐藏变量的边际分布和后验分布。通过消息传递，在树的叶子节点到根节点的方向上收集证据，反向通过根节点到叶子节点分发信息。
- 变量消元 (Variable Elimination, VE)**
  - 关键词：**符号化计算、填充边、消元顺序、和弦图、最大团、边际计算
  - 用途：**VE用于在图模型中计算边际和条件概率。通过消去变量并逐步计算因子，减少计算复杂度。通过添加填充边，将图转化为和弦图，使得最大团的计算更加高效。
- 联接树算法 (Junction Tree Algorithm, JTA)**
  - 关键词：**和弦图、联接树、最大团、RIP（运行交集属性）、树上的推断
  - 用途：**JTA是信念传播的推广，适用于复杂图结构。通过构建和弦图、提取最大团并构建联接树，可以进行推断。它适用于一般图（不仅是树形结构），并解决了VE在复杂图中的推断问题。
- 推断的计算复杂度**
  - 关键词：**计算不可行性、NP难、#P难、3-SAT问题、精确推断

- **用途：**推断在图模型中的计算复杂度非常高，尤其是在涉及**树宽度**较大的图中。  
精确推断是**NP难的**，且**#P难**，表明对于某些模型，计算量极其庞大，不适合使用精确推断。

## 5. 近似推断方法

- **关键词：**启发式方法、近似推断、无保证的准确性、快速推断
  - **用途：**对于那些不能使用精确推断的图模型，近似推断方法可以提供有效的解决方案，尽管没有准确性保证。这些方法通常用于具有复杂结构的模型或计算不可行的情况。
- 

## ② 子章节关联

### 1. 20.3 精确推断的计算复杂性 与 20.4 联接树算法：

- **关系：**20.3节讲述了变量消元 (VE) 和联接树算法 (JTA) 的计算复杂度，表明这些算法在某些图结构（如树宽度较大的图）下可能变得计算不可行。因此，20.4节通过介绍联接树算法作为推断的一种有效方法，解决了复杂图模型中的推断问题。

### 2. 20.4 联接树算法 与 20.5 近似推断：

- **关系：**20.4节介绍了联接树算法，这是一种用于复杂图模型的推断方法。然而，尽管JTA在某些模型中表现良好，但对于一些大规模或复杂的图模型，JTA仍然可能不可行。因此，20.5节转向介绍近似推断方法，这些方法不保证准确性，但能在许多实际问题中提供良好的解决方案。

### 3. 20.4 联接树算法 与 20.5.1 近似推断：

- **关系：**在20.4节中，介绍了JTA作为图模型推断的有效方法，但其计算复杂度仍然可能很高。20.5.1节介绍了对于精确推断不可行的情形，**近似推断方法**成为一种解决方案。
- 

## ③ 本章总结

本章的重点是讨论图模型中的**推断算法**及其计算复杂性。具体分为三大部分：

1. **精确推断：**介绍了\*\*信念传播 (BP) 和变量消元 (VE) \*\*算法，重点是如何通过这些方法在树形和链形结构的图中计算边际概率和后验概率。尽管这些算法在某些模型中非常有效，但它们的计算复杂度可能非常高，尤其是在处理复杂图结构时。
2. **联接树算法 (JTA)：**针对复杂图模型，JTA作为一种推广了信念传播的算法，通过将图转换为和弦图、提取最大团并构建联接树来提高推断效率。尽管如此，JTA在最坏情况下的计算复杂度仍然可能非常大，尤其是当图的**树宽度**很高时。
3. **计算不可行性与近似推断：**对于计算复杂度过高的推断任务，近似推断方法提供了一种有效的解决方案。虽然这些方法没有精确度和时间保证，但在很多实际问题中，它们通常能提供良好的结果。

总体来说，本章展示了图模型中推断的多种方法，从精确推断到近似推断，同时强调了在不同图结构下的计算复杂度问题和现实应用中的挑战。

请用通俗易懂的话总结本章，使我能够理解

本章的核心内容是关于**图模型推断**的不同方法，重点讨论了推断的**计算难度**以及如何使用不同的算法来解决这些问题。下面是通俗易懂的总结：

### 1. 精确推断的挑战：

- 图模型是一种用来表示变量及其相互关系的数学模型，推断就是通过已知的信息来推测未知的变量。
- \*\*信念传播 (BP) 和变量消元 (VE)\*\* 是两种经典的推断方法，能够帮助我们计算图模型中隐藏变量的概率。
- 但是，这些方法在处理**复杂图结构**时非常耗时，尤其是当图的**树宽度**较大时，计算时间可能会呈指数增长，变得几乎不可行。

### 2. 联接树算法 (JTA) :

- 为了解决复杂图模型中的推断问题，我们可以使用**联接树算法 (JTA)**。它是一种优化版的信念传播算法，能够在复杂图中更有效地进行推断。
- JTA通过将图转换为**和弦图**并构建**联接树**来简化推断过程。尽管如此，在最坏的情况下，JTA的计算复杂度依然可能非常高。

### 3. 计算不可行性与近似推断：

- 即使使用JTA，某些图模型的推断依然是不可行的。这时候我们需要使用**近似推断方法**，这些方法通过折中处理来快速得到一个近似解，虽然它们不能保证完全准确。
- 近似推断方法虽然在理论上没有精确推断那么强的保证，但在实际应用中，它们往往能够提供足够好的解，并且运行得更快。

### 总结：

本章讲解了在图模型中进行推断时可能遇到的计算困难，并提供了**几种推断方法**。对于一些简单的图模型，信念传播和变量消元可以提供精确的结果。但对于复杂的图，推断变得非常困难，我们需要依赖像**联接树算法**和**近似推断方法**这样的优化算法，尽管这些方法在最坏情况下仍然会遇到计算瓶颈。