

Course notes for MIT EECS 6.819/6.869 Spring 2022

Antonio Torralba, Phillip Isola and William Freeman

April 18, 2022

Contents

28 Probabilistic graphical models	5
28.1 Simple Examples	6
28.2 Directed graphical models	9
29 Inference in graphical models	11
29.1 Simple example of inference in a graphical model	12
29.2 Belief propagation (BP)	13
29.2.1 Derivation of message-passing rule	13
29.2.2 Marginal probability	14
29.2.3 Message update sequence	14
29.2.4 Example belief propagation application: stereo	18
29.3 Loopy belief propagation	21
29.4 MAP estimation and energy models	21
29.4.1 Numerical example of Belief Propagation	22

Chapter 28

Probabilistic graphical models

Probabilistic graphical models describe joint probability distributions in a modular way that allows us to reason about the visual world even when we're modeling very complicated situations. These models are useful in vision, where we often need to exploit modularity to make computations tractable.

A probabilistic graphical model is a graph that describes a class of probability distributions sharing a common structure. The graph has nodes, drawn as circles, indicating the variables of the joint probability. It has edges, drawn as lines connecting nodes to other nodes. At first, we'll restrict our attention to a type of graphical model called undirected, so the edges are line segments without arrowheads.

The edges indicate the conditional independence structure of the nodes of the graph. If there is no edge between two nodes, then the variables described by those nodes are independent, if you condition on the values of intervening nodes in any path between them in the graph. If two nodes do have a line between them, then we cannot assume they are independent. We introduce this through a set of examples.

28.1 Simple Examples

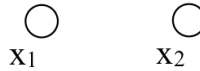
Here is the simplest probabilistic graphical model:



(28.1)

This “graph” is just a single node for the variable x_1 . There is no restriction imposed by the graph on the functional form of its probability function. In keeping with notation we’ll use later, we write that the probability distribution over x_1 , $P(x_1) = \phi_1(x_1)$.

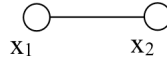
Another trivial graphical model is this:



(28.2)

Here, the lack of a line connecting x_1 with x_2 indicates a lack of statistical dependence between these two variables. In detail, we condition on knowing the values of any connecting neighbors of x_1 and x_2 —in this case there are none—and thus x_1 and x_2 are statistically independent. Because of that independence, the joint probability depicted by this graphical model must be a product of each variable’s marginal probability and thus must have this form: $P(x_1, x_2) = \phi_1(x_1)\phi_2(x_2)$.

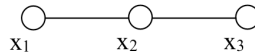
Let’s add a line between these two variables:



(28.3)

By that graph, we mean that there may be a statistical dependency between the two variables. The class of probability distributions depicted here is now more general than the one above, and we can only write that the joint probability is some unknown function of x_1 and x_2 , $P(x_1, x_2) = \phi_{12}(x_1, x_2)$. This graph offers no simplification from the most general probability function of two variables.

Here is a graph with some structure:

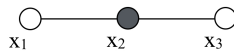


(28.4)

This graph means that if we condition on the variable x_2 , then x_1 and x_3 are independent. A general form for $P(x_1, x_2, x_3)$ that guarantees such conditional independence is

$$P(x_1, x_2, x_3) = \phi_{12}(x_1, x_2)\phi_{23}(x_2, x_3). \quad (28.5)$$

Note the conditional independence implicit in Eq. (28.5): if the value of x_2 is given (indicated by the solid circle in the graph below), then the joint probability of Eq. (28.5) becomes a product of some function $f(x_1)$ times some other function $g(x_3)$, revealing the conditional independence of x_1 and x_3 . Graphically, we denote the conditioning on the variable x_2 with a filled circle for that node:



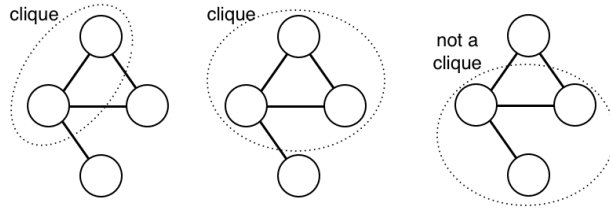
(28.6)

In the next chapter, we’ll exploit that structure of the joint probability to perform inference efficiently using belief propagation.

A celebrated theorem, the **Hammersley-Clifford theorem** [Besag1974], tells the form the joint probability must have for any given graphical model. The joint probability for a

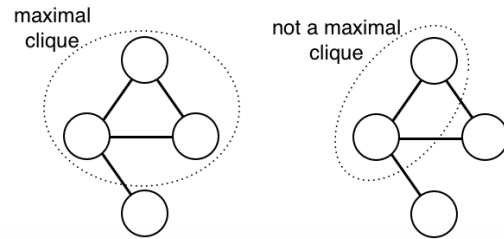
probabilistic graphical model must be a product of functions of each the “maximal cliques” of the graph. Here we define the terms.

A **clique** is any set of nodes where each node is connected to every other node in the clique. These graphs illustrate the clique property:



(28.7)

A maximal clique is a clique that can't include more nodes of the graph without losing the clique property. The sets of nodes below form maximal cliques (left), or do not (right):



(28.8)

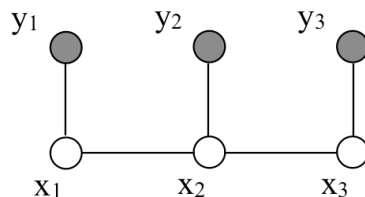
The **Hammersley-Clifford theorem** [Hammersley and Clifford1971]: A positive probability distribution has the independence structure described by a graphical model if and only if it can be written as a product of functions over the variables of each maximal clique:

$$P(x_1, x_2, \dots, x_N) = \prod_{x_c \in x_i} \Psi_c(x_c), \quad (28.9)$$

where the product is over all maximal cliques x_c in the graph, x_i .

In the example of the graph 28.4, the maximal cliques are (x_1, x_2) and (x_2, x_3) , so the Hammersley-Clifford theorem says that the corresponding joint probability must be of the form of Eq. (28.5).

Now we examine some graphical model structures that are especially useful in vision. In perception problems, we typically have both observed and unobserved variables. The graph below shows a simple “Markov chain” [Gagnieu2017] structure with 3 observed variables, shaded, labeled y_i , and 3 unobserved variables, labeled x_i :



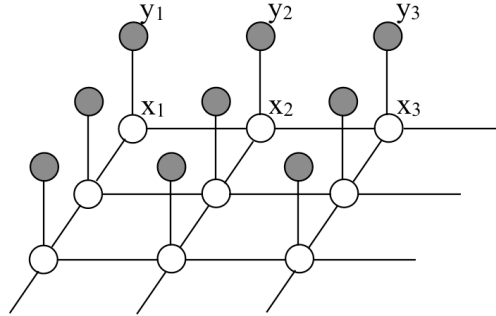
(28.10)

This is a chain because the variables form a linear sequence. It's Markov because the hidden variables have the Markov property: conditioned on x_2 , variable x_3 is independent of variable x_1 . The joint probability of all the variables shown here is $P(x_1, x_2, x_3, y_1, y_2, y_3) = \phi_{12}(x_1, x_2)\phi_{23}(x_2, x_3)\psi_1(y_1, x_1)\psi_2(y_2, x_2)\psi_3(y_3, x_3)$. Using $P(a, b) = P(a|b)P(b)$, we can also write the probability of the x variables conditioned on the observations y ,

$$P(x_1, x_2, x_3|y_1, y_2, y_3) = \frac{1}{P(\mathbf{y})} \phi_{12}(x_1, x_2)\phi_{23}(x_2, x_3)\psi_1(y_1, x_1)\psi_2(y_2, x_2)\psi_3(y_3, x_3). \quad (28.11)$$

(For brevity, we write $P(\mathbf{y})$ for $P(y_1, y_2, y_3)$.) Thus, to form the conditional distribution, within a normalization factor, we simply include the observed variable values into the joint probability. For vision applications, we often use such Markov chain structures to describe events over time.

To capture relationships over space, a 2-d structure is useful, called a **Markov random field**, or MRF [A Blake2011]:



(28.12)

Then the joint probability over all the variables factorizes into this product:

$$P(\mathbf{x}|\mathbf{y}) = \frac{1}{P(\mathbf{y})} \prod_{(i,j)} \phi_{ij}(x_i, x_j) \prod_i \psi_i(x_i, y_i), \quad (28.13)$$

where the first product is over all spatial neighbors i and j , and the second product is over all nodes i .

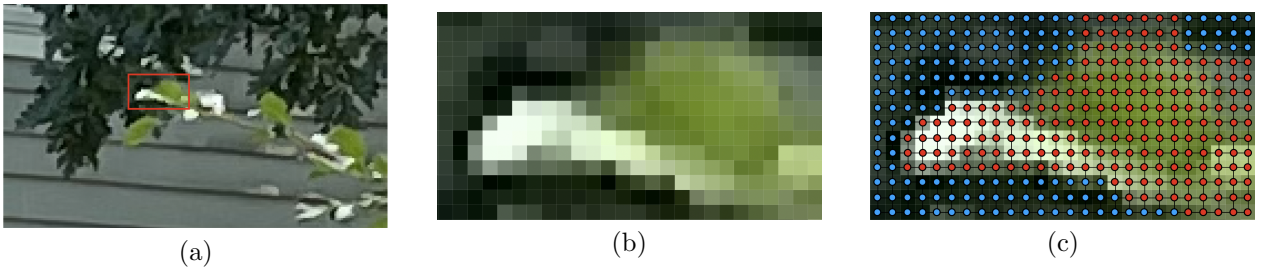
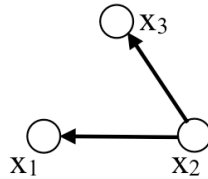


Figure 28.1: (a) A local region of an image to be segmented (b) Enlarged region (c) Visualization of a MRF of nodes corresponding to image pixels. The internal states of the nodes are indicator variables of segment membership. A hypothetical most probable configuration of the indicator variables are shown as the filled color of each node's circle.

An example of how Markov random field models are applied to images is depicted in Fig. 28.1. (a) and (b) show a small image region and its context in the image. (c) shows an MRF with each node corresponding to a pixel of the image region. The states of an MRF can be indicator variables of image segment membership for each pixel. The states of a hypothetical most-probable configuration of this MRF are shown as the color of each node in this illustration.

28.2 Directed graphical models

In addition to undirected graphical models, another type of graphical model is commonly used. **Directed graphical models** [Koller and Friedman2009] describe factorizations of the joint probability into products of conditional probability distributions. Each node in a directed graph contributes a well-specified factor in the joint probability: the probability of its variable, conditioned all the variables originating arrows pointing into it. So this graph:



(28.14)

denotes the joint probability,

$$P(x_1, x_2, x_3) = P(x_2)P(x_1|x_2)P(x_3|x_2) \quad (28.15)$$

The general rule for writing the joint probability described by a directed graph is this: Each node, x_n , contributes the factor $P(x_n|x_{\Xi})$, where Ξ is the set of nodes with arrows pointing in to node x_n . You can verify that Eq. 28.15 follows this rule. Directed graphical models are often used to describe causal processes.

Chapter 29

Inference in graphical models

Given a probabilistic graphical model and observations, we want to estimate the states of the unobserved variables. For example, given image observations, we want to estimate the pose of the person. The belief propagation algorithm lets us do that efficiently.

Recall the Bayesian inference task: our observations are the elements of a vector, \mathbf{y} , and we seek to infer the probability $P(\mathbf{x}|\mathbf{y})$ of some scene parameters, \mathbf{x} , given the observations. By Bayes' rule, we have

$$P(\mathbf{x}|\mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{x})P(\mathbf{x})}{P(\mathbf{y})} \quad (29.1)$$

The *likelihood term* $P(\mathbf{y}|\mathbf{x})$ describes how a rendered scene \mathbf{x} generates observations \mathbf{y} . The *prior probability*, $P(\mathbf{x})$, tells the probability of any given scene \mathbf{x} occurring. For inference, we often ignore the denominator, $P(\mathbf{y})$, called the *evidence*, as it is constant with respect to the variables we seek to estimate, \mathbf{x} .

Typically, we make many observations, \mathbf{y} , of the variables of some system, and we want to find the the state of some hidden variable, \mathbf{x} , given those observations. The posterior probability, $P(\mathbf{x}|\mathbf{y})$, tells us the probability for any value of the hidden variables, \mathbf{x} . From this posterior probability, we often want some single *best* estimate for \mathbf{x} , denoted $\hat{\mathbf{x}}$ and called a point estimate.

Selecting the best estimate $\hat{\mathbf{x}}$ requires specifying a penalty for making a wrong guess. If we penalize all wrong answers equally, the best strategy is to guess the value of \mathbf{x} that maximizes the posterior probability, $P(\mathbf{x}|\mathbf{y})$ (because any other explanation $\hat{\mathbf{x}}$ for the observations \mathbf{y} would be less probable). That is called the MAP estimate, for *maximum a posteriori*.

But we may want to penalize wrong answers as a function of how far they are from the correct answer. If that penalty function is the squared distance in \mathbf{x} , then the point estimate that minimizes the average value of that error is called the minimum mean squared error estimate, or MMSE. To find this estimate, we seek the $\hat{\mathbf{x}}$ that minimizes the squared error, weighted by the probability of each outcome:

$$\hat{\mathbf{x}}_{MMSE} = \operatorname{argmin}_{\tilde{\mathbf{x}}} \int_{\mathbf{x}} P(\mathbf{x}|\mathbf{y})(\mathbf{x} - \tilde{\mathbf{x}})'(\mathbf{x} - \tilde{\mathbf{x}})d\mathbf{x} \quad (29.2)$$

Differentiating with respect to \mathbf{x} to solve for the stationary point, the global minimum for this convex function, we find

$$\hat{\mathbf{x}}_{MMSE} = \int_{\mathbf{x}} \mathbf{x}P(\mathbf{x}|\mathbf{y})d\mathbf{x}. \quad (29.3)$$

Thus, the minimum mean square error estimate, $\hat{\mathbf{x}}_{MMSE}$, is the mean of the posterior distribution. If \mathbf{x} represents a discretized space, then the marginalization integrals over $d\mathbf{x}$ become sums over the corresponding discrete states of \mathbf{x} .

For now, we'll assume we seek the MMSE estimate. By the properties of the multi-variate mean, to find the mean at each variable, or node, in a network, we can first find the marginal

probability at each node, then compute the mean of each marginal probability. In other words, given $P(\mathbf{x}|\mathbf{y})$, we will compute $P(x_i|\mathbf{y})$, where i is the i th hidden node. From $P(x_i|\mathbf{y})$ it is simple to compute the posterior mean at node i .

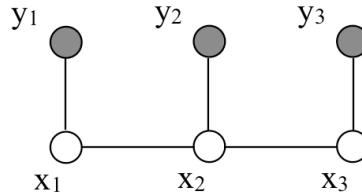
For the case of discrete variables, we compute the marginal probability at a node by summing over the states at all the other nodes,

$$P(x_i|\mathbf{y}) = \sum_{j \setminus i} \sum_{x_j} P(x_1, x_2, \dots, x_i, \dots, x_N | \mathbf{y}), \quad (29.4)$$

where the notation $j \setminus i$ means "all possible values of j except for $j = i$ ".

29.1 Simple example of inference in a graphical model

To gain intuition, let's calculate the marginal probability for a simple example. Consider the three-node Markov chain of 28.10, reproduced here.



(29.5)

Vision tasks this can apply to include modeling the probability of a pixel belonging to an object edge, given the evidence, observations \mathbf{y} , at a point and two neighboring locations. The inferred states \mathbf{x} could be a label indicating the presence of an edge.

We seek to marginalize the joint probability in order to find the marginal probability at node 1, $P(x_1|\mathbf{y})$, given observations y_1, y_2 , and y_3 , which we denote as \mathbf{y} . We have, assuming the nodes have discrete states,

$$P(x_1|\mathbf{y}) = \sum_{x_2} \sum_{x_3} P(x_1, x_2, x_3 | \mathbf{y}) \quad (29.6)$$

Here's the main point: If we knew nothing about the structure of the joint probability $P(x_1, x_2, x_3 | \mathbf{y})$, the computation would require $|x|^3$ computations: a double-sum over all x_2 and x_3 states must be computed for each possible x_1 state. (We're denoting the number of states of any of the x variables as $|x|$). In the more general case, for a Markov chain of N nodes, we would need $|x|^N$ summations to compute the desired marginal at any node of the chain. Such a computation quickly becomes intractable as N grows.

But we can exploit the model's structure to avoid the exponential growth of the computation with N . Substituting the joint probability, from the graphical model, into the marginalization equation, Eq. (29.6), gives

$$P(x_1|\mathbf{y}) = \frac{1}{P(\mathbf{y})} \sum_{x_2} \sum_{x_3} \phi_{12}(x_1, x_2) \phi_{23}(x_2, x_3) \psi_1(y_1, x_1) \psi_2(y_2, x_2) \psi_3(y_3, x_3) \quad (29.7)$$

This form for the joint probability reveals that not every variable is coupled to every other one. We can pass summations through variables they don't sum over, letting us compute the marginalization much more efficiently. This will make only a small difference for this short chain, but it makes a huge difference for longer ones. So we write

$$P(x_1|\mathbf{y}) = \frac{1}{P(\mathbf{y})} \sum_{x_2} \sum_{x_3} \phi_{12}(x_1, x_2) \phi_{23}(x_2, x_3) \psi_1(y_1, x_1) \psi_2(y_2, x_2) \psi_3(y_3, x_3) \quad (29.8)$$

$$= \frac{1}{P(\mathbf{y})} \psi_1(y_1, x_1) \sum_{x_2} \phi_{12}(x_1, x_2) \psi_2(y_2, x_2) \sum_{x_3} \phi_{23}(x_2, x_3) \psi_3(y_3, x_3) \quad (29.9)$$

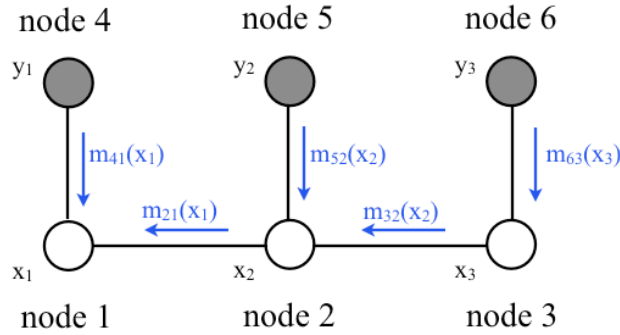


Figure 29.1: Summary of the messages (partial sums) for a simple belief propagation example.

That factorization of Eq. (29.9) is the key step: It reduces the number of terms summed from order $|x|^3$ to order $2|x|^2$ for this chain, and for a chain of length N chain, from order $|x|^N$ to order $(N-1)|x|^2$, a huge computational savings for large N .

The partial sums of Eq. (29.9) are named “messages” because they pass information from one node to another. We call the message from node 3 to node 2, $m_{32}(x_2) = \sum_{x_3} \phi_{23}(x_2, x_3) m_{63}(x_3)$. The other partial sum is the message from node 2 to node 1, $m_{21}(x_1) = \sum_{x_2} \phi_{12}(x_1, x_2) m_{52}(x_2) m_{32}(x_2)$. Note that the messages are always messages about the states of the node that the message is being sent *to*—the arguments of the message m_{ij} are the states x_j of node j . The algorithm corresponding to Eq. (29.9) is called *belief propagation*.

Eq. (29.9) gives us the marginal probability at node 1. To find the marginal probability at another node, you can write out the sums over variables needed for that node, pass the sums through factors in the joint probability that they don’t operate on, to come up with an efficient reorganization of summations analogous to Eq. (29.9). You would find that many of the summations from marginalization at the node x_1 would need to be recomputed for the marginalization at node x_2 . That motivates storing and reusing the messages, which the belief propagation algorithm does in an optimal way.

29.2 Belief propagation (BP)

For more complicated graphical models, we want to replace the manual factorization with an automatic procedure for identifying the computations needed for marginalizations and to cache them efficiently. Belief propagation does that by identifying those reusable sums, the “messages”.

29.2.1 Derivation of message-passing rule

We’ll describe belief propagation (BP) only for the special case of graphical models with pair-wise potentials. The clique potentials between neighboring nodes are $\psi_{ij}(x_j, x_i)$. Extensions to higher-order potentials is straightforward. (Convert the graphical model into one with only pairwise potentials. This can be done by augmenting the state of some nodes to encompass several nodes, until the remaining nodes only need pairwise potential functions in their factorization of the joint probability.) You can find formal derivations of belief propagation in [Jordan1998, Koller and Friedman2009].

Consider Fig. 29.2, showing a section of a general network with pair-wise potentials. There is a network of $N+1$ nodes, numbered 0 through N and we will marginalize over nodes $x_1 \dots x_N$. Fig. 29.2 (a) shows the marginalization equation for a network of variables with discrete states. (For continuous variables, integrals replace the summations). If we assume the nodes form a tree, we can distribute the marginalization sum past nodes for which the sum is a constant value to obtain the sums depicted in Fig. 29.2 (b).

We define a *belief propagation message*:

A message m_{ij} , from node i to node j , is the sum of the probability over all states of all nodes in the subtree that leaves node i and does not include node j .

Referring to Fig. 29.2, (a) shows the desired marginalization sum over a network of pairwise cliques. (This approach can be extended to more general cliques). We can pass the summations over node states through nodes that are constant over those summations, arriving at the factorization shown in (b). Remembering that the message from node i to node j is the sum over the tree leaving node i , we can then read-off the recursive belief propagation message update rule by inspecting Fig. 29.2:

To compute the message from node j to node i :

1. Multiply together all messages coming in to node j , except for the message from node i back to node j ,
2. Multiply by the pairwise compatibility function $\psi_{ij}(x_i, x_j)$.
3. Marginalize over the variable x_j .

These steps are summarized in this equation to compute the message from node j to node i :

$$m_{ji}(x_i) = \sum_{x_j} \psi_{ij}(x_i, x_j) \prod_{k \in \eta(j) \setminus i} m_{kj}(x_j) \quad (29.10)$$

where $\eta(j) \setminus i$ means “the neighbors of node j except for node i ”. Figure 29.3 shows this equation in a graphical form. Any local potential functions $\phi_j(x_j)$ are treated as an additional message into node j , $m_{0j}(x_j) = \phi_j(x_j)$. For the case of continuous variables, the sum over the states of x_j in Eq. (29.10) is replaced by an integral over the domain of x_j .

As mentioned above, BP messages are partial sums in the marginalization calculation. The arguments of messages are always the state of the node that the message is going to. (Belief propagation follows the “nosey neighbor rule” (from Brendan Frey, U. Toronto). Every node is a house in some neighborhood. Your nosey neighbor says to you: “given what I’ve seen myself and everything I’ve heard from my neighbors, here’s what I think is going on inside your house.” That metaphor made sense to me after my children became teenagers.)

29.2.2 Marginal probability

The marginal probability at a node i , Eq. (29.4), is the sum over the joint probabilities of all states except those of x_i . Because we assume the network has no loops, the conditional independence structure assures us that marginal probability at a node i is the product of the sum of all states of all nodes in each sub-tree connected to node i . (Conditioned on node i , the probabilities within each subtree are independent). Thus the marginal probability at node i is the product of all the incoming messages:

$$P_i(x_i) = \prod_{j \in \eta(i)} m_{ji}(x_i) \quad (29.11)$$

(We include the local clique potential at node i , $\psi_i(x_i)$, as one of the messages in the product of all messages into node i).

29.2.3 Message update sequence

To find all the messages, how do we invoke the recursive BP update rule, Eq. (29.10)? We can apply Eq. (29.10) whenever all the incoming messages in the BP update rule are defined. If there are no incoming messages to a node, then its outgoing message is well-defined in the update rule. This lets us start the recursive algorithm.

A node can send a message whenever all the incoming messages it needs have been computed. We can compute the outgoing messages from leaf nodes in the graphical model

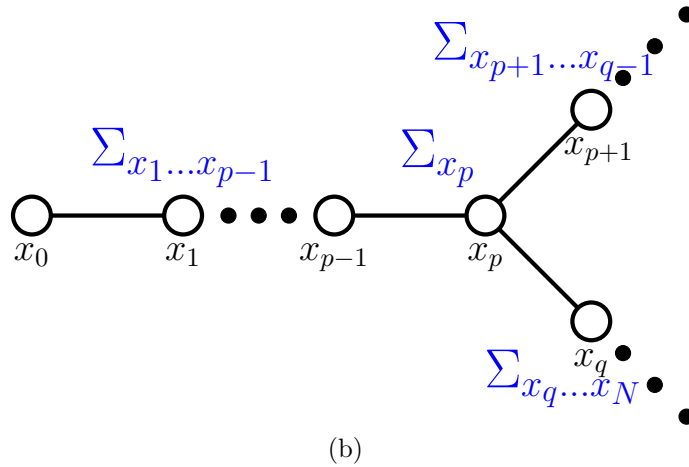
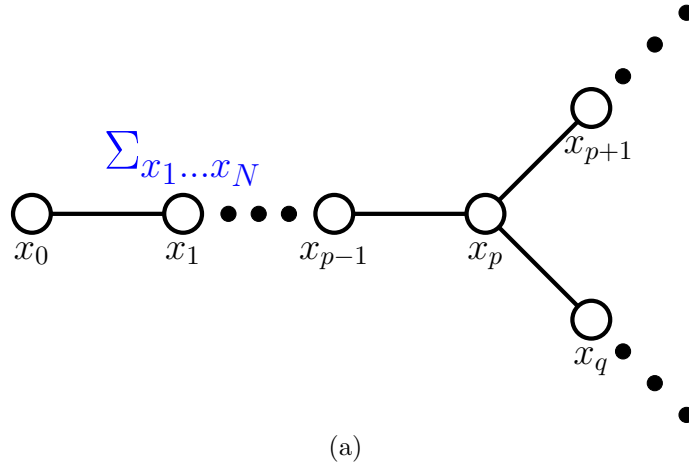


Figure 29.2: Generic example motivating the formula for belief propagation. (a) shows the marginalization for a general graph with no loops, focussing on the partial sums at node x_p . (b) Shows how the partial sums distribute over the nodes. Note that, by the definition of a message, $m_{p+1,p}$ and m_{qp} correspond to these partial sums over nodes indicated in the graph: $m_{p+1,p} = \sum_{x_{p+1} \dots x_{q-1}}$ and $m_{qp} = \sum_{x_q \dots x_N}$. Marginalization over the states of all nodes leaving node p , not including node $p-1$, leads to Eq. (29.10) for the belief propagation message passing rule.

$$\begin{array}{c} m_{ji}(x_i) \\ \left| \right. \\ \left| \right. \\ \left| \right. \end{array} = \begin{array}{c} \Phi_{ij}(x_i, x_j) \\ \left[\begin{array}{c} x_j \\ x_i \end{array} \right] \\ \times \end{array} \begin{array}{c} m_{k_1j}(x_j) \\ \left| \right. \\ \left| \right. \\ \left| \right. \end{array} \cdot \begin{array}{c} m_{k_2j}(x_j) \\ \left| \right. \\ \left| \right. \\ \left| \right. \end{array} \cdot \begin{array}{c} m_{k_3j}(x_j) \\ \left| \right. \\ \left| \right. \\ \left| \right. \end{array} \cdot \dots$$

Figure 29.3: Pictorial depiction of belief propagation message passing rules of Eq. (29.10), showing the vector and matrix shapes. To send a message from node j to node i : We term-by-term multiply (shown by \cdot) the messages (column vectors) coming in to node j , then matrix multiply (shown by \times) the resulting column vector by the compatibility matrix $\Phi_{ij}(x_i, x_j)$ to obtain $m_{ji}(x_i)$.

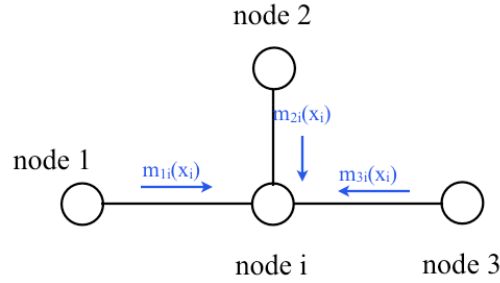


Figure 29.4: To compute the marginal probability at node i , we multiply together all the incoming messages at that node: $P_i(x_i) = \prod_{j \in \eta(i)} m_{ji}(x_i)$, remembering to include any local potential terms $\phi_i(x_i)$ as another message.

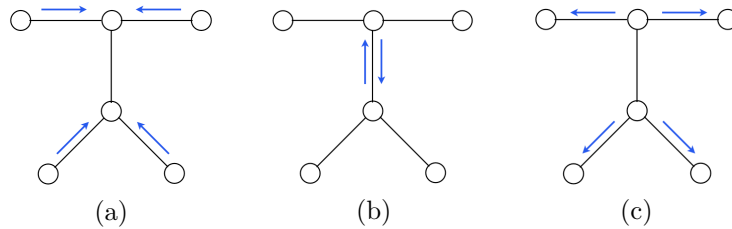


Figure 29.5: Example of a **synchronous parallel update schedule** for BP message passing. Whenever any node has the required incoming messages needed to send an outgoing message, it does. (a) At the first iteration, only the leaf nodes have the needed incoming messages to send an outgoing message (by definition, leaf nodes have no links other than the one on which they'll be sending their outgoing message, so they have no incoming messages to wait for). (b) second iteration, (c) third iteration. By the third iteration, every edge has messages computed in both directions, and we can now compute the marginal probability at every node in the graph.

tree, since they have no incoming messages other than from the node to which they are sending a message, which doesn't enter in the outgoing message computation. Two natural message passing protocols are consistent with that rule: depth-first update, and parallel update. In depth-first update, one node is arbitrarily picked as the root. Messages are then passed from the leaves of the tree (leaves, relative to that root node) up to the root, then back down to the leaves. In parallel update, at each turn, every node sends every outgoing message for which it has received all the necessary incoming messages. Figure 29.5 depicts the flow of messages for the parallel, synchronous update scheme.

Note that, when computing marginals at many nodes, we re-use messages with the BP algorithm. A single message-passing sweep through all the nodes lets us calculate the marginal at any node (using the depth-first update rules to calculate the marginal at the root node). A second sweep from the root node back to all the leaf nodes calculates all the messages needed to find the marginal probability at every node. It takes only twice the number of computations to calculate the incoming messages to every node as it does to calculate the incoming messages at a single node.

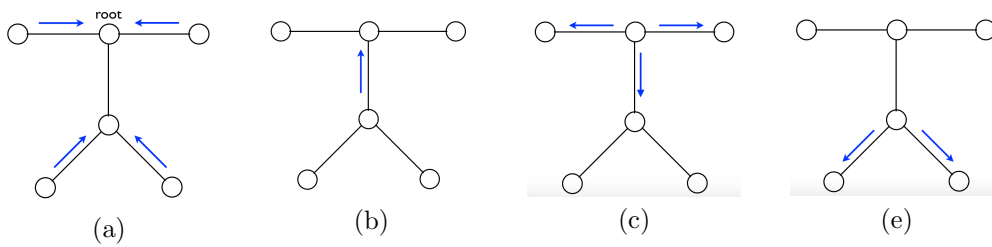


Figure 29.6: Example of a **depth-first update schedule** for BP message passing. We arbitrarily select one node as the root. (a) Messages begin at the leaves and (b) proceed to the root. Once all the messages reach the designated root node, (c) an outgoing sweep computes the remaining messages, (d) ending at the leaf nodes.

29.2.4 Example belief propagation application: stereo

Assume that we want to compute the depth from a stereo image pair, such as the pair shown in (a) and (b) in Fig. 29.7, and with their insets marked with the red rectangles, shown in (c) and (d). Assume that the two images have been rectified, as described in Chapter (stereo), and consider the common scanline, marked in each of (c) and (d) with a black horizontal line. The luminance intensities of each scanline are plotted in Fig. 29.8. Most computer vision stereo algorithms [Scharstein and Szeliski2002] examine multiple scanlines to compute the depth at every pixel, but to illustrate a simple example of belief propagation in vision, we will consider the stereo vision problem while considering only one scanline pair at a time.

The graphical model that describes this stereo depth reconstruction for a single scanline is shown in Fig. 29.9. Each unfilled circle in the graphical models represents the unknown depth of each pixel in the scan line of one of the stereo images, in this example, the left camera image in Fig. 29.7, (c). The open circles mean that the depth of each pixels is an unobserved variable.

The intensity values of the left and right camera scan lines are used to compute the local evidence for the disparity offset, $d[i, j]$, of each right camera pixel corresponding to the left image pixel at each right camera position, $[i, j]$. For this example, we assume that each right camera pixel value, x_r is that of a corresponding left camera value, x_l , but with independent, identically distributed (IID) Gaussian random noise added. This leads to a simple formula for the local evidence, ψ_i , for the given depth value of any pixel,

$$\psi_i = P(d[i, j]|x_r, x_l) = ke^{-\sum_{i=-10}^{i=10} \frac{(x_r[i, j] - x_l[i - d[i, j], j])^2}{2\sigma^2}} \quad (29.12)$$

For simplicity in this example, we discretize the x_r to x_l disparity offset into 4 different bins, each 45 disparity offsets wide. For simplicity, we add the local disparity evidence over the 45 depth states, then normalize the sum at each n to maintain a probability distribution. The resulting local evidence matrix, for each of the 4 depth states, at each left camera spatial position, is displayed in the 3rd row of Fig. 29.10. The intensities in each column in the third-row image add-up to one.

The prior probability of any configuration of pixel disparity states is described by setting the compatibility matrices, $\phi[s_i, s_{i+1}]$ of the Markov chain of Fig. 29.9. For this example, we use a compatibility matrix referred to as the Potts Model [Wu1982]:

$$\phi[s_i, s_{i+1}] = \begin{cases} 1, & \text{if } s_i = s_{i+1} \\ \delta, & \text{otherwise} \end{cases} \quad (29.13)$$

For this example, we used $\delta = 0.001$

We then ran BP, Eq. 29.10 in a depth-first update schedule. For this network, that involves two linear sweeps over all the nodes of the chain, first from left to right, then from right to left. Starting from the left-most node, we updated the rightward messages one node at a time, in a rightward sweep, then, starting from the right-most node, updated all the leftward messages one node at a time, in a leftward sweep.

Once all the messages were computed, we used Eq. 29.11 to find the posterior marginal probability at each node. The local evidence, the leftward and rightward messages, and the final marginal probability at each node are displayed for every left-camera pixel of the inset. The four rows of each display in subimages (c)-(f) correspond to the four possible depth states used for this calculation. Note that for this scanline, in this image pair, there is a depth discontinuity at location of the red vertical mark in Fig. 29.10—the canoe is closer to the camera than the ground beyond the canoe that appears to the right of the canoe. The local evidence for depth disparity, (c), is rather inaccurate, but the marginal posterior probability accurately finds the depth discontinuity at the canoe boundary.

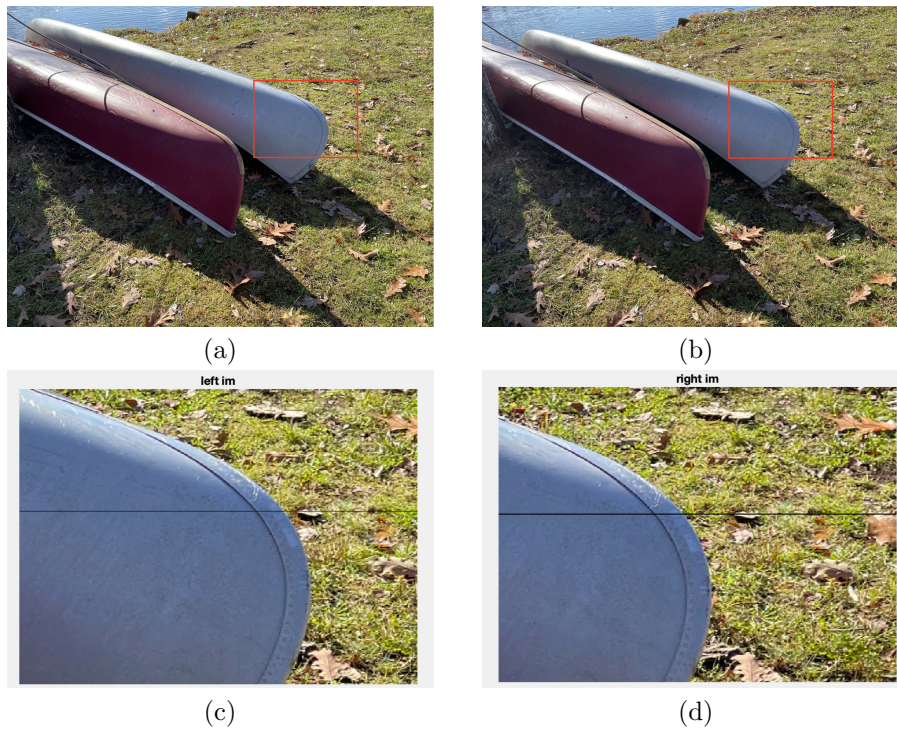


Figure 29.7: (a) left camera image and (b) right camera image. (c) and (d): insets showing areas of analysis.

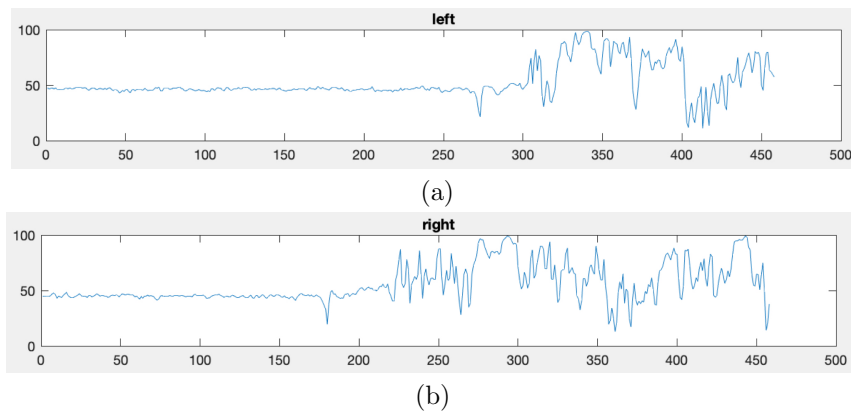


Figure 29.8: (a) and (b) show the luminance traces from the scanlines marked in black in Fig. 29.7 (c) and (d), respectively.

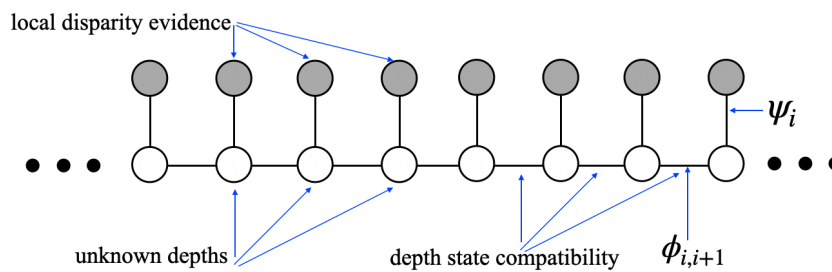


Figure 29.9: Graphical model for the posterior probability for stereo disparity offset between the left and right camera views from a stereo rig.

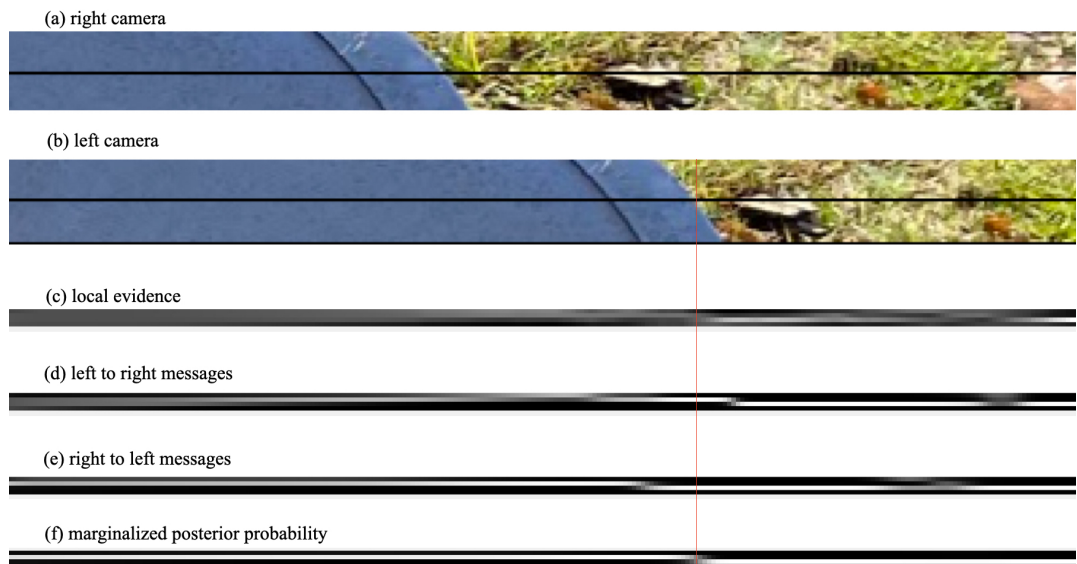


Figure 29.10: Belief propagation applied to graphical model, Fig. 29.9, for the stereo problem. (a) and (b): left and right camera views. The luminance values covered by the black line in each figure is the stereo observation for this example. (c) is the average over a depth range of the local evidence for a given depth offset disparity. (d) and (e) show the local messages computed in the belief propagation algorithm at each position. (f) shows the final marginalized posterior probability at each pixel of the left camera. Note that the belief propagation algorithm accurately finds the position of the depth discontinuity.

29.3 Loopy belief propagation

The BP message update rules only work to give the exact marginals when the topology of the network is that of a tree or a chain. In general, one can show that exact computation of marginal probabilities for graphs with loops depends on a graph-theoretic quantity known as the treewidth of the graph. For many graphical models of interest in vision, such as 2-d Markov Random Fields related to images, these quantities can be intractably large.

But the message update rules are described locally, and one might imagine that it is a useful local operation to perform, even without the global guarantees of ordinary BP. It turns out that is true. Here is the algorithm: **loopy belief propagation algorithm**

1. convert graph to pairwise potentials
2. initialize all messages to all ones, or to random values between 0 and 1.
3. run the belief propagation update rules of Sect. 29.2.1 until convergence.

One can show that, in networks with loops, fixed points of the belief propagation algorithm (message configurations where the messages don't change with a message update) correspond to minima of a well-known approximation from the statistical physics community known as the Bethe free energy [Yedidia et al.2001]. In practice, the solutions found by the loopy belief propagation algorithm are often quite good [Kevin Murphy1999].

29.4 MAP estimation and energy models

Instead of summing over the states of other nodes, we are sometimes interested in finding the \mathbf{x} that maximizes the joint probability. The argmax operator passes through constant variables just as the summation sign did. This leads to an alternate version of the belief propagation algorithm, Eq. (29.10), with the summation (of multiplying the vector message products by the node compatibility matrix) replaced with “argmax”. This is called the “max-product” version of belief propagation, and it computes an MAP estimate of the hidden states.

Improvements have been developed over loopy belief propagation for the case of MAP estimation, see, for example, tree-reweighted belief propagation [Kolmogorov2006] and graph cuts [Zabih and Komogorov2004].

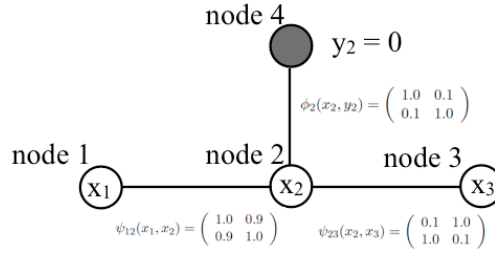


Figure 29.11: Undirected graphical model used in belief propagation example

29.4.1 Numerical example of Belief Propagation

Here we work through a numerical example of belief propagation. To make the arithmetic easy, we'll solve for the marginal probabilities in the graphical model of two-state (0 and 1) random variables shown in Fig. 29.11. That graphical model has 3 hidden variables, and one variable observed to be in state 0. The compatibility matrices are given in the arrays below (for which the state indices are 0, then 1, reading from left to right and top to bottom).

$$\psi_{12}(x_1, x_2) = \begin{pmatrix} 1.0 & 0.9 \\ 0.9 & 1.0 \end{pmatrix} \quad (29.14)$$

$$\psi_{23}(x_2, x_3) = \begin{pmatrix} 0.1 & 1.0 \\ 1.0 & 0.1 \end{pmatrix} \quad (29.15)$$

$$\psi_{42}(x_2, y_2) = \begin{pmatrix} 1.0 & 0.1 \\ 0.1 & 1.0 \end{pmatrix} \quad (29.16)$$

Note that in defining these potential functions, we haven't taken care to normalize the joint probability, so we'll need to normalize each marginal probability at the end. (remember $P(x_1, x_2, x_3, y_2) = \psi_{42}(x_2, y_2)\psi_{23}(x_2, x_3)\psi_{12}(x_1, x_2)$, which should sum to 1 after summing over all states.)

For this simple toy example, we can tell what results to expect by inspection, then verify that BP is doing the right thing. Node x_2 wants very much to look like $y_2 = 0$, because $\psi_{42}(x_2, y_2)$ contributes a large valued to the posterior probability when $x_2 = y_2 = 1$ or when $x_2 = y_2 = 0$. From $\psi_{12}(x_1, x_2)$ we see that x_1 has a very mild preference to look like x_2 . So we expect the marginal probability at node x_2 will be heavily biased toward $x_2 = 0$, and that node x_1 will have a mild preference for state 0. $\psi_{23}(x_2, x_3)$ encourages x_3 to be the opposite of x_2 , so it will be biased toward the state $x_3 = 1$.

Let's see what belief propagation gives. We'll follow the parallel, synchronous update scheme for calculating all the messages. The leaf nodes can send messages along their edges without waiting for any messages to be updated. For the message from node 1, we have

$$m_{12}(x_2) = \sum_{x_1} \psi_{12}(x_1, x_2) \quad (29.17)$$

$$= \begin{pmatrix} 1.0 & 0.9 \\ 0.9 & 1.0 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad (29.18)$$

$$= \begin{pmatrix} 1.9 \\ 1.9 \end{pmatrix} \quad (29.19)$$

$$= k \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad (29.20)$$

For numerical stability, we typically normalize the computed messages in Eq. (29.10) so the entries sum to 1, or so their maximum entry is 1, then remember to renormalize the final marginal probabilities to sum to 1. Here, we've normalized the messages for simplicity, (absorbing the normalization into a constant, k).

The message from node 3 to node 2 is

$$m_{32}(x_2) = \sum_{x_3} \psi_{32}(x_2, x_3) \quad (29.21)$$

$$= \begin{pmatrix} 0.1 & 1.0 \\ 1.0 & 0.1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad (29.22)$$

$$= \begin{pmatrix} 1.1 \\ 1.1 \end{pmatrix} \quad (29.23)$$

$$= k \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad (29.24)$$

We have a non-trivial message from observed node y_2 (node 4) to the hidden variable x_2 :

$$m_{42}(x_2) = \sum_{x_4} \psi_{42}(x_2, y_2) \quad (29.25)$$

$$= \begin{pmatrix} 1.0 & 0.1 \\ 0.1 & 1.0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad (29.26)$$

$$= \begin{pmatrix} 1.0 \\ 0.1 \end{pmatrix} \quad (29.27)$$

where y_2 has been fixed to $y_2 = 0$, thus restricting $\psi_{42}(x_2, y_2)$ to just the first column.

Now we just have two messages left to compute before we have all messages computed (and therefore all node marginals computed from simple combinations of those messages). The message from node 2 to node 1 uses the messages from nodes 4 to 2 and 3 to 2:

$$m_{21}(x_1) = \sum_{x_2} \psi_{12}(x_1, x_2) m_{42}(x_2) m_{32}(x_2) \quad (29.28)$$

$$= \begin{pmatrix} 1.0 & 0.9 \\ 0.9 & 1.0 \end{pmatrix} \left[\begin{pmatrix} 1.0 \\ 0.1 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right] = \begin{pmatrix} 1.09 \\ 1.0 \end{pmatrix} \quad (29.29)$$

The final message is that from node 2 to node 3 (since y_2 is observed, we don't need to compute the message from node 2 to node 4). That message is:

$$m_{23}(x_3) = \sum_{x_2} \psi_{23}(x_2, x_3) m_{42}(x_2) m_{12}(x_2) \quad (29.30)$$

$$= \begin{pmatrix} 0.1 & 1.0 \\ 1.0 & 0.1 \end{pmatrix} \left[\begin{pmatrix} 1.0 \\ 0.1 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right] = \begin{pmatrix} 0.2 \\ 1.01 \end{pmatrix} \quad (29.31)$$

Now that we've computed all the messages, let's look at the marginals of the three hidden nodes. The product of all the messages arriving at node 1 is just the one message, $m_{21}(x_1)$, so we have (introducing constant k to normalize the product of messages to be a probability distribution)

$$P_1(x_1) = k m_{21}(x_1) = \frac{1}{2.09} \begin{pmatrix} 1.09 \\ 1.0 \end{pmatrix} \quad (29.32)$$

As we knew it should, node 1 shows a slight preference for state 0.

The marginal at node 2 is proportional to the product of 3 messages. Two of those are trivial messages, but we'll show them all for completeness:

$$P_2(x_2) = k m_{12}(x_2) m_{42}(x_2) m_{32}(x_2) \quad (29.33)$$

$$= k \begin{pmatrix} 1 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} 1.0 \\ 0.1 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad (29.34)$$

$$= \frac{1}{1.1} \begin{pmatrix} 1.0 \\ 0.1 \end{pmatrix} \quad (29.35)$$

As expected, belief propagation reveals a strong bias for node 2 being in state 0.

Finally, for the marginal probability at node 3, we have

$$P_3(x_3) = km_{23}(x_3) = \frac{1}{1.21} \begin{pmatrix} 0.2 \\ 1.01 \end{pmatrix} \quad (29.36)$$

As predicted, this variable is biased toward being in state 1.

By running belief propagation within this tree, we have computed the exact marginal probabilities at each node, reusing the intermediate sums across different marginalizations, and exploiting the structure of the joint probability to perform the computation efficiently. If nothing were known about the joint probability structure, the marginalization cost would grow exponentially with the number of nodes in the network. But if the graph structure corresponding to the joint probability is known to be a chain or a tree, then the marginalization cost only grows linearly with the number of nodes, and is quadratic in the node state dimensions. The belief propagation algorithm enables inference in many large-scale problems.

Bibliography

- [A Blake2011] A Blake P Kohli, C Rother. 2011. *Markov random fields for vision and image processing*. MIT Press.
- [Besag1974] Besag, Julian. 1974. Spatial interaction and the statistical analysis of lattice systems. *Journal of the Royal Statistical Society. Series B (Methodological)* 36 (2): 192–236. <http://www.jstor.org/stable/2984812>.
- [Gagniuc2017] Gagniuc, P. A. 2017. *Markov chains: From theory to implementation and experimentation*. John Wiley Sons..
- [Hammersley and Clifford1971] Hammersley, J. M., and P. Clifford. 1971. Markov fields on finite graphs and lattices. <http://www.statslab.cam.ac.uk/~grg/books/hammfest/hamm-cliff.pdf>.
- [Jordan1998] Jordan, M. I., ed. 1998. *Learning in graphical models*. MIT Press.
- [Kevin Murphy1999] Kevin Murphy Yair Weiss, Michael I. Jordan. 1999. Loopy belief propagation for approximate inference: An empirical study. In *Proceedings of the fifteenth conference on uncertainty in artificial intelligence (uai1999)*. <https://arxiv.org/abs/1301.6725>.
- [Koller and Friedman2009] Koller, D., and N. Friedman, eds. 2009. *Probabilistic graphical models*. MIT Press.
- [Kolmogorov2006] Kolmogorov, V. 2006. Convergent tree-reweighted message passing for energy minimization. *IEEE Pattern Analysis and Machine Intelligence* 28 (10).
- [Scharstein and Szeliski2002] Scharstein, D., and R. Szeliski. 2002. A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *International Journal of Computer Vision* 47: 7–42. <https://vision.middlebury.edu/stereo/>.
- [Wu1982] Wu, Fa-Yueh. 1982. The potts model. *Rev. Mod. Phys.* 54 (1): 235–268.
- [Yedidia et al.2001] Yedidia, J. S., W. T. Freeman, and Y. Weiss. 2001. Generalized belief propagation, Vol. 13, 689–695. MIT Press.
- [Zabih and Komogorov2004] Zabih, R., and V. Komogorov. 2004. What energy functions can be minimized via graph cuts? In *European conf. computer vision*, Vol. 26, 147–159.