

Variational Message Passing

John Winn

Christopher M. Bishop

Microsoft Research Cambridge

Roger Needham Building

7 J. J. Thomson Avenue

Cambridge CB3 0FB, U.K.

JWINN@MICROSOFT.COM

CMBISHOP@MICROSOFT.COM

Editor: Tommi Jaakkola

Abstract

Bayesian inference is now widely established as one of the principal foundations for machine learning. In practice, exact inference is rarely possible, and so a variety of approximation techniques have been developed, one of the most widely used being a deterministic framework called variational inference. In this paper we introduce Variational Message Passing (VMP), a general purpose algorithm for applying variational inference to Bayesian Networks. Like belief propagation, VMP proceeds by sending messages between nodes in the network and updating posterior beliefs using local operations at each node. Each such update increases a lower bound on the log evidence (unless already at a local maximum). In contrast to belief propagation, VMP can be applied to a very general class of conjugate-exponential models because it uses a factorised variational approximation. Furthermore, by introducing additional variational parameters, VMP can be applied to models containing non-conjugate distributions. The VMP framework also allows the lower bound to be evaluated, and this can be used both for model comparison and for detection of convergence. Variational message passing has been implemented in the form of a general purpose inference engine called VIBES ('Variational Inference for BayEasian networkS') which allows models to be specified graphically and then solved variationally without recourse to coding.

Keywords: Bayesian networks, variational inference, message passing

1. Introduction

Variational inference methods (Neal and Hinton, 1998; Jordan et al., 1998) have been used successfully for a wide range of models, and new applications are constantly being explored. In each previous application, the equations for optimising the variational approximation have been worked out by hand, a process which is both time consuming and error prone. For several other inference methods, general purpose algorithms have been developed which can be applied to large classes of probabilistic models. For example, *belief propagation* can be applied to any acyclic discrete network (Pearl, 1986) or mixed-Gaussian network (Lauritzen, 1992), and the Monte Carlo algorithm described in Thomas et al. (1992) can perform Gibbs sampling in almost any Bayesian network. Similarly, *expectation propagation* (Minka, 2001) has been successfully applied to a wide range of models. Each of these algorithms relies on being able to decompose the required computation into calculations that are local to each node in the graph and which require only messages passed along the edges connected to that node.

However, Monte Carlo methods are computationally very intensive, and also suffer from difficulties in diagnosing convergence, while belief propagation is only guaranteed to converge for tree-structured graphs. Expectation propagation is limited to certain classes of model for which the required expectations can be evaluated, is also not guaranteed to converge in general, and is prone to finding poor solutions in the case of multi-modal distributions. For these reasons the framework of variational inference has received much attention.

In this paper, the variational message passing algorithm is developed, which optimises a variational bound using a set of local computations for each node, together with a mechanism for passing messages between the nodes. VMP allows variational inference to be applied automatically to a large class of Bayesian networks, without the need to derive application-specific update equations. In VMP, the messages are exponential family distributions, summarised either by their natural parameter vector (for child-to-parent messages) or by a vector of moments (for parent-to-child messages). These messages are defined so that the optimal variational distribution for a node can be found by summing the messages from its children together with a function of the messages from its parents, where this function depends on the conditional distribution for the node.

The VMP framework applies to models described by directed acyclic graphs in which the conditional distributions at each node are members of the exponential family, which therefore includes discrete, Gaussian, Poisson, gamma, and many other common distributions as special cases. For example, VMP can handle a general DAG of discrete nodes, or of linear-Gaussian nodes, or an arbitrary combination of the two provided there are no links going from continuous to discrete nodes. This last restriction can be lifted by introducing further variational bounds, as we shall discuss. Furthermore, the marginal distribution of observed data represented by the graph is not restricted to be from the exponential family, but can come from a very broad class of distributions built up from exponential family building blocks. The framework therefore includes many well known machine learning algorithms such as hidden Markov models, probabilistic PCA, factor analysis and Kalman filters, as well as mixtures and hierarchical mixtures of these.

Note that, since we work in a fully Bayesian framework, latent variables and parameters appear on an equal footing (they are all unobserved stochastic variables which are marginalised out to make predictions). If desired, however, point estimates for parameters can be made simply by maximising the same bound as used for variational inference.

As an illustration of the power of the message passing viewpoint, we use VMP within a software tool called VIBES (Variational Inference in BayEsian networkS) which allows a model to be specified by drawing its graph using a graphical interface, and which then performs variational inference automatically on this graph.

The paper is organised as follows. Section 2 gives a brief review of variational inference methods. Section 3 contains the derivation of the variational message passing algorithm, along with an example of its use. In Section 4 the class of models which can be handled by the algorithm is defined, while Section 5 describes the VIBES software. Some extensions to the algorithm are given in Section 6, and Section 7 concludes with an overall discussion and suggestions for future research directions.

2. Variational Inference

In this section, variational inference will be reviewed briefly with particular focus on the case where the variational distribution has a factorised form. The random variables in the model will be denoted

by $\mathbf{X} = (\mathbf{V}, \mathbf{H})$ where \mathbf{V} are the visible (observed) variables and \mathbf{H} are the hidden (latent) variables. We assume that the model has the form of a Bayesian network and so the joint distribution $P(\mathbf{X})$ can be expressed in terms of the conditional distributions at each node i ,

$$P(\mathbf{X}) = \prod_i P(X_i | \text{pa}_i) \tag{1}$$

where pa_i denotes the set of variables corresponding to the parents of node i and X_i denotes the variable or group of variables associated with node i .

Ideally, we would like to perform exact inference within this model to find posterior marginal distributions over individual latent variables. Unfortunately, exact inference algorithms, such as the junction tree algorithm (Cowell et al., 1999), are typically only applied to discrete or linear-Gaussian models and are computationally intractable for all but the simplest models. Instead, we must turn to approximate inference methods. Here, we consider the deterministic approximation method of variational inference.

The goal in variational inference is to find a tractable variational distribution $Q(\mathbf{H})$ that closely approximates the true posterior distribution $P(\mathbf{H} | \mathbf{V})$. To do this we note the following decomposition of the log marginal probability of the observed data, which holds for any choice of distribution: $Q(\mathbf{H})$

$$\ln P(\mathbf{V}) = \mathcal{L}(Q) + \text{KL}(Q || P). \tag{2}$$

Here

$$\begin{aligned} \mathcal{L}(Q) &= \sum_{\mathbf{H}} Q(\mathbf{H}) \ln \frac{P(\mathbf{H}, \mathbf{V})}{Q(\mathbf{H})} \\ \text{KL}(Q || P) &= - \sum_{\mathbf{H}} Q(\mathbf{H}) \ln \frac{P(\mathbf{H} | \mathbf{V})}{Q(\mathbf{H})} \end{aligned} \tag{3}$$

and the sums are replaced by integrals in the case of continuous variables. Here $\text{KL}(Q || P)$ is the Kullback-Leibler divergence between the true posterior $P(\mathbf{H} | \mathbf{V})$ and the variational approximation $Q(\mathbf{H})$. Since this satisfies $\text{KL}(Q || P) \geq 0$ it follows from (2) that the quantity $\mathcal{L}(Q)$ forms a lower bound on $\ln P(\mathbf{V})$.

We now choose some family of distributions to represent $Q(\mathbf{H})$ and then seek a member of that family that maximises the lower bound $\mathcal{L}(Q)$ and hence minimises the Kullback-Leibler divergence $\text{KL}(Q || P)$. If we allow $Q(\mathbf{H})$ to have complete flexibility then we see that the maximum of the lower bound occurs when the Kullback-Leibler divergence is zero. In this case, the variational posterior distribution equals the true posterior and $\mathcal{L}(Q) = \ln P(\mathbf{V})$. However, working with the true posterior distribution is computationally intractable (otherwise we wouldn't be resorting to variational methods). We must therefore consider a more restricted family of Q distributions which has the property that the lower bound (3) can be evaluated and optimised efficiently and yet which is still sufficiently flexible as to give a good approximation to the true posterior distribution.

This method leads to minimisation of 'exclusive' divergence $\text{KL}(Q || P)$ rather than the 'inclusive' divergence $\text{KL}(P || Q)$. Minimising the exclusive divergence can lead to a Q which ignores modes of P . However, minimising the inclusive divergence can lead to Q assigning posterior mass to regions where P has vanishing density. If the latter behaviour is preferred, then there are alternative approximation techniques for minimising the inclusive divergence, including expectation propagation (Minka, 2001).

2.1 Factorised Variational Distributions

We wish to choose a variational distribution $Q(\mathbf{H})$ with a simpler structure than that of the model, so as to make calculation of the lower bound $\mathcal{L}(Q)$ tractable. One way to simplify the dependency structure is by choosing a variational distribution where disjoint groups of variables are independent. This is equivalent to choosing Q to have a factorised form

$$Q(\mathbf{H}) = \prod_i Q_i(\mathbf{H}_i), \quad (4)$$

where $\{\mathbf{H}_i\}$ are the disjoint groups of variables. This approximation has been successfully used in many applications of variational methods (Attias, 2000; Ghahramani and Beal, 2001; Bishop, 1999). Substituting (4) into (3) gives

$$\mathcal{L}(Q) = \sum_{\mathbf{H}} \prod_i Q_i(\mathbf{H}_i) \ln P(\mathbf{H}, \mathbf{V}) - \sum_i \sum_{\mathbf{H}_i} Q_i(\mathbf{H}_i) \ln Q_i(\mathbf{H}_i).$$

We now separate out all terms in one factor Q_j ,

$$\begin{aligned} \mathcal{L}(Q) &= \sum_{\mathbf{H}_j} Q_j(\mathbf{H}_j) \langle \ln P(\mathbf{H}, \mathbf{V}) \rangle_{\sim Q_j(\mathbf{H}_j)} + \mathbb{H}(Q_j) + \sum_{i \neq j} \mathbb{H}(Q_i) \\ &= -\text{KL}(Q_j \| Q_j^*) + \text{terms not in } Q_j \end{aligned} \quad (5)$$

where \mathbb{H} denotes entropy and we have introduced a new distribution Q_j^* , defined by

$$\ln Q_j^*(\mathbf{H}_j) = \langle \ln P(\mathbf{H}, \mathbf{V}) \rangle_{\sim Q(\mathbf{H}_j)} + \text{const.} \quad (6)$$

and $\langle \cdot \rangle_{\sim Q(\mathbf{H}_j)}$ denotes an expectation with respect to all factors except $Q_j(\mathbf{H}_j)$. The bound is maximised with respect to Q_j when the KL divergence in (5) is zero, which occurs when $Q_j = Q_j^*$. Therefore, we can maximise the bound by setting Q_j equal to Q_j^* . Taking exponentials of both sides we obtain

$$Q_j^*(\mathbf{H}_j) = \frac{1}{Z} \exp \langle \ln P(\mathbf{H}, \mathbf{V}) \rangle_{\sim Q(\mathbf{H}_j)}, \quad (7)$$

where Z is the normalisation factor needed to make Q_j^* a valid probability distribution. Note that the equations for all of the factors are coupled since the solution for each $Q_j(\mathbf{H}_j)$ depends on expectations with respect to the other factors $Q_{i \neq j}$. The variational optimisation proceeds by initialising each of the $Q_j(\mathbf{H}_j)$ and then cycling through each factor in turn replacing the current distribution with a revised estimate given by (7).

3. Variational Message Passing

In this section, the variational message passing algorithm will be derived and shown to optimise a factorised variational distribution using a message passing procedure on a graphical model. For the initial derivation, it will be assumed that the variational distribution is factorised with respect to each hidden variable H_i and so can be written

$$Q(\mathbf{H}) = \prod_i Q_i(H_i).$$

From (6), the optimised form of the j th factor is given by

$$\ln Q_j^*(H_j) = \langle \ln P(\mathbf{H}, \mathbf{V}) \rangle_{\sim Q(H_j)} + \text{const.}$$

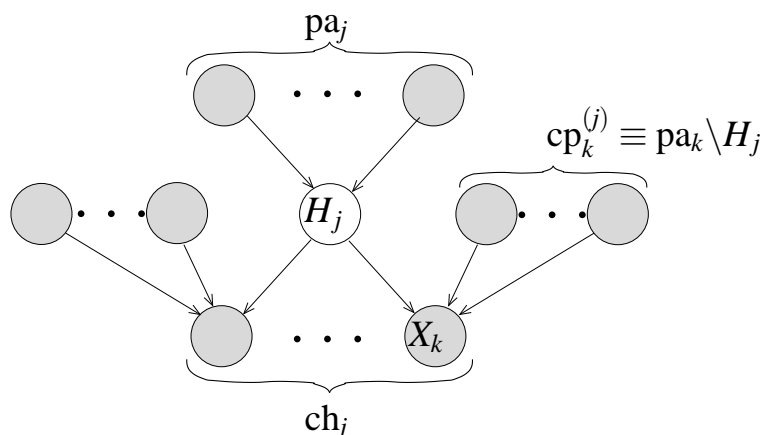


Figure 1: A key observation is that the variational update equation for a node H_j depends only on expectations over variables in the Markov blanket of that node (shown shaded), defined as the set of parents, children and co-parents of that node.

We now substitute in the form of the joint probability distribution of a Bayesian network, as given in (1),

$$\ln Q_j^*(H_j) = \langle \sum_i \ln P(X_i | \text{pa}_i) \rangle_{\sim Q(H_j)} + \text{const.}$$

Any terms in the sum over i that do not depend on H_j will be constant under the expectation and can be subsumed into the constant term. This leaves only the conditional $P(H_j | \text{pa}_j)$ together with the conditionals for all the children of H_j , as these have H_j in their parent set,

$$\ln Q_j^*(H_j) = \langle \ln P(H_j | \text{pa}_j) \rangle_{\sim Q(H_j)} + \sum_{k \in \text{ch}_j} \langle \ln P(X_k | \text{pa}_k) \rangle_{\sim Q(H_j)} + \text{const.} \quad (8)$$

where ch_j are the children of node j in the graph. Thus, the expectations required to evaluate Q_j^* involve only those variables lying in the Markov blanket of H_j , consisting of its parents, children and co-parents¹ $\text{cp}_k^{(j)}$. This is illustrated in the form of a directed graphical model in Figure 1. Note that we use the notation X_k to denote both a random variable and the corresponding node in the graph. The optimisation of Q_j can therefore be expressed as a local computation at the node H_j . This computation involves the sum of a term involving the parent nodes, along with one term from each of the child nodes. These terms can be thought of as ‘messages’ from the corresponding nodes. Hence, we can decompose the overall optimisation into a set of local computations that depend only on messages from neighbouring (i.e. parent and child) nodes in the graph.

3.1 Conjugate-Exponential Models

The exact form of the messages in (8) will depend on the functional form of the conditional distributions in the model. It has been noted (Attias, 2000; Ghahramani and Beal, 2001) that important simplifications to the variational update equations occur when the distributions of variables, condi-

1. The co-parents of a node X are all the nodes with at least one child which is also a child of X (excluding X itself).

tioned on their parents, are drawn from the exponential family and are conjugate² with respect to the distributions over these parent variables. A model where both of these constraints hold is known as a *conjugate-exponential* model.

A conditional distribution is in the exponential family if it can be written in the form

$$P(\mathbf{X} | \mathbf{Y}) = \exp[\phi(\mathbf{Y})^T \mathbf{u}(\mathbf{X}) + f(\mathbf{X}) + g(\mathbf{Y})] \quad (9)$$

where $\phi(\mathbf{Y})$ is called the *natural parameter* vector and $\mathbf{u}(\mathbf{X})$ is called the *natural statistic* vector. The term $g(\mathbf{Y})$ acts as a normalisation function that ensures the distribution integrates to unity for any given setting of the parameters \mathbf{Y} . The exponential family contains many common distributions, including the Gaussian, gamma, Dirichlet, Poisson and discrete distributions. The advantages of exponential family distributions are that expectations of their logarithms are tractable to compute and their state can be summarised completely by the natural parameter vector. The use of conjugate distributions means that the posterior for each factor has the same form as the prior and so learning changes only the values of the parameters, rather than the functional form of the distribution.

If we know the natural parameter vector $\phi(\mathbf{Y})$ for an exponential family distribution, then we can find the expectation of the natural statistic vector with respect to the distribution. Rewriting (9) and defining \tilde{g} as a reparameterisation of g in terms of ϕ gives,

$$P(\mathbf{X} | \phi) = \exp[\phi^T \mathbf{u}(\mathbf{X}) + f(\mathbf{X}) + \tilde{g}(\phi)].$$

We integrate with respect to \mathbf{X} ,

$$\int_{\mathbf{X}} \exp[\phi^T \mathbf{u}(\mathbf{X}) + f(\mathbf{X}) + \tilde{g}(\phi)] d\mathbf{X} = \int_{\mathbf{X}} P(\mathbf{X} | \phi) d\mathbf{X} = 1$$

and then differentiate with respect to ϕ

$$\begin{aligned} \int_{\mathbf{X}} \frac{d}{d\phi} \exp[\phi^T \mathbf{u}(\mathbf{X}) + f(\mathbf{X}) + \tilde{g}(\phi)] d\mathbf{X} &= \frac{d}{d\phi}(1) = 0 \\ \int_{\mathbf{X}} P(\mathbf{X} | \phi) \left[\mathbf{u}(\mathbf{X}) + \frac{d\tilde{g}(\phi)}{d\phi} \right] d\mathbf{X} &= 0. \end{aligned}$$

And so the expectation of the natural statistic vector is given by

$$\langle \mathbf{u}(\mathbf{X}) \rangle_{P(\mathbf{X} | \phi)} = -\frac{d\tilde{g}(\phi)}{d\phi}. \quad (10)$$

We will see later that the factors of our Q distribution will also be in the exponential family and will have the same natural statistic vector as the corresponding factor of P . Hence, the expectation of \mathbf{u} under the Q distribution can also be found using (10).

3.2 Optimisation of Q in Conjugate-Exponential Models

We will now demonstrate how the optimisation of the variational distribution can be carried out, given that the model is conjugate-exponential. We consider the general case of optimising a factor

2. A parent distribution $P(X | Y)$ is said to be *conjugate* to a child distribution $P(W | X)$ if $P(X | Y)$ has the same functional form, with respect to X , as $P(W | X)$.

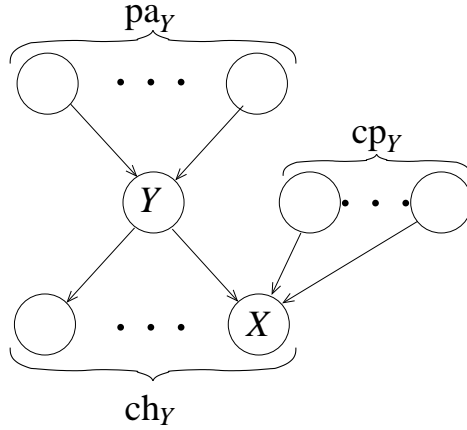


Figure 2: Part of a graphical model showing a node Y , the parents and children of Y , and the co-parents of Y with respect to a child node X .

$Q(Y)$ corresponding to a node Y , whose children include X , as illustrated in Figure 2. From (9), the log conditional probability of the variable Y given its parents can be written

$$\ln P(Y | \text{pa}_Y) = \phi_Y(\text{pa}_Y)^T \mathbf{u}_Y(Y) + f_Y(Y) + g_Y(\text{pa}_Y). \quad (11)$$

The subscript Y on each of the functions $\phi_Y, \mathbf{u}_Y, f_Y, g_Y$ is required as these functions differ for different members of the exponential family and so need to be defined separately for each node.

Consider a node $X \in \text{ch}_Y$ which is a child of Y . The conditional probability of X given its parents will also be in the exponential family and so can be written in the form

$$\ln P(X | Y, \text{cp}_Y) = \phi_X(Y, \text{cp}_Y)^T \mathbf{u}_X(X) + f_X(X) + g_X(Y, \text{cp}_Y) \quad (12)$$

where cp_Y are the co-parents of Y with respect to X , in other words, the set of parents of X excluding Y itself. The quantity $P(Y | \text{pa}_Y)$ in (11) can be thought of as a prior over Y , and $P(X | Y, \text{cp}_Y)$ as a (contribution to) the likelihood of Y .

Example

If X is Gaussian distributed with mean Y and precision β , it follows that the co-parent set cp_Y contains only β , and the log conditional for X is

$$\ln P(X | Y, \beta) = \begin{bmatrix} \beta Y \\ -\beta/2 \end{bmatrix}^T \begin{bmatrix} X \\ X^2 \end{bmatrix} + \frac{1}{2}(\ln \beta - \beta Y^2 - \ln 2\pi). \quad (13)$$

Conjugacy requires that the conditionals of (11) and (12) have the same functional form with respect to Y , and so the latter can be rewritten in terms of $\mathbf{u}_Y(Y)$ by defining functions ϕ_{XY} and λ as follows

$$\ln P(X | Y, \text{cp}_Y) = \phi_{XY}(X, \text{cp}_Y)^T \mathbf{u}_Y(Y) + \lambda(X, \text{cp}_Y). \quad (14)$$

It may appear from this expression that the function ϕ_{XY} depends on the form of the parent conditional $P(Y | \text{pa}_Y)$ and so cannot be determined locally at X . This is not the case, because the conjugacy constraint dictates $\mathbf{u}_Y(Y)$ for any parent Y of X , implying that ϕ_{XY} can be found directly from the form of the conditional $P(X | \text{pa}_X)$.

Example Continuing the above example, we can find ϕ_{XY} by rewriting the log conditional in terms of Y to give

$$\ln P(X|Y, \beta) = \begin{bmatrix} \beta X \\ -\beta/2 \end{bmatrix}^T \begin{bmatrix} Y \\ Y^2 \end{bmatrix} + \frac{1}{2}(\ln \beta - \beta X^2 - \ln 2\pi),$$

which lets us define ϕ_{XY} and dictate what $\mathbf{u}_Y(Y)$ must be to enforce conjugacy,

$$\phi_{XY}(X, \beta) \stackrel{\text{def}}{=} \begin{bmatrix} \beta X \\ -\beta/2 \end{bmatrix}, \quad \mathbf{u}_Y(Y) = \begin{bmatrix} Y \\ Y^2 \end{bmatrix}. \quad (15)$$

From (12) and (14), it can be seen that $\ln P(X|Y, \text{cp}_Y)$ is linear in $\mathbf{u}_X(X)$ and $\mathbf{u}_Y(Y)$ respectively. Conjugacy also dictates that this log conditional will be linear in $\mathbf{u}_Z(Z)$ for each co-parent $Z \in \text{cp}_Y$. Hence, $\ln P(X|Y, \text{cp}_Y)$ must be a multi-linear³ function of the natural statistic functions \mathbf{u} of X and its parents. This result is general, for any variable A in a conjugate-exponential model, the log conditional $\ln P(A|\text{pa}_A)$ must be a multi-linear function of the natural statistic functions of A and its parents.

Example The log conditional $\ln P(X|Y, \beta)$ in (13) is multi-linear in each of the vectors,

$$\mathbf{u}_X(X) = \begin{bmatrix} X \\ X^2 \end{bmatrix}, \quad \mathbf{u}_Y(Y) = \begin{bmatrix} Y \\ Y^2 \end{bmatrix}, \quad \mathbf{u}_\beta(\beta) = \begin{bmatrix} \beta \\ \ln \beta \end{bmatrix}.$$

Returning to the variational update equation (8) for a node Y , it follows that all the expectations on the right hand side can be calculated in terms of the $\langle \mathbf{u} \rangle$ for each node in the Markov blanket of Y . Substituting for these expectations, we get

$$\begin{aligned} \ln Q_Y^*(Y) &= \langle \phi_Y(\text{pa}_Y)^T \mathbf{u}_Y(Y) + f_Y(Y) + g_Y(\text{pa}_Y) \rangle_{\sim Q(Y)} \\ &+ \sum_{k \in \text{ch}_Y} \langle \phi_{XY}(X_k, \text{cp}_k)^T \mathbf{u}_Y(Y) + \lambda(X_k, \text{cp}_k) \rangle_{\sim Q(Y)} + \text{const.} \end{aligned}$$

which can be rearranged to give

$$\begin{aligned} \ln Q_Y^*(Y) &= \left[\langle \phi_Y(\text{pa}_Y) \rangle_{\sim Q(Y)} + \sum_{k \in \text{ch}_Y} \langle \phi_{XY}(X_k, \text{cp}_k) \rangle_{\sim Q(Y)} \right]^T \mathbf{u}_Y(Y) \\ &+ f_Y(Y) + \text{const.} \end{aligned} \quad (16)$$

It follows that Q_Y^* is an exponential family distribution of the same form as $P(Y|\text{pa}_Y)$ but with a natural parameter vector ϕ_Y^* such that

$$\phi_Y^* = \langle \phi_Y(\text{pa}_Y) \rangle + \sum_{k \in \text{ch}_Y} \langle \phi_{XY}(X_k, \text{cp}_k) \rangle \quad (17)$$

where all expectations are with respect to Q . As explained above, the expectations of ϕ_Y and each ϕ_{XY} are multi-linear functions of the expectations of the natural statistic vectors corresponding to their dependent variables. It is therefore possible to reparameterise these functions in terms of these

3. A function f is a multi-linear function of variables a, b, \dots if it varies linearly with respect to each variable, for example, $f(a, b) = ab + 3b$ is multi-linear in a and b . Although, strictly, this function is *affine* in a because of the constant term, we follow common usage and refer to it as linear.

expectations

$$\begin{aligned}\tilde{\phi}_Y(\{\langle \mathbf{u}_i \rangle\}_{i \in \text{pa}_Y}) &= \langle \phi_Y(\text{pa}_Y) \rangle \\ \tilde{\phi}_{XY}(\langle \mathbf{u}_k \rangle, \{\langle \mathbf{u}_j \rangle\}_{j \in \text{cp}_k}) &= \langle \phi_{XY}(X_k, \text{cp}_k) \rangle.\end{aligned}$$

The final step is to show that we can compute the expectations of the natural statistic vectors \mathbf{u} under Q . From (16) any variable A has a factor Q_A with the same exponential family form as $P(A | \text{pa}_A)$. Hence, the expectations of \mathbf{u}_A can be found from the natural parameter vector of that distribution using (10). In the case where A is observed, the expectation is irrelevant and we can simply calculate $\mathbf{u}_A(A)$ directly.

Example In (15), we defined $\phi_{XY}(X, \beta) = \begin{bmatrix} \beta X \\ -\beta/2 \end{bmatrix}$. We now reparameterise it as

$$\tilde{\phi}_{XY}(\langle \mathbf{u}_X \rangle, \langle \mathbf{u}_\beta \rangle) \stackrel{\text{def}}{=} \begin{bmatrix} \langle \mathbf{u}_\beta \rangle_0 \langle \mathbf{u}_X \rangle_0 \\ -\frac{1}{2} \langle \mathbf{u}_\beta \rangle_0 \end{bmatrix}$$

where $\langle \mathbf{u}_X \rangle_0$ and $\langle \mathbf{u}_\beta \rangle_0$ are the first elements of the vectors $\langle \mathbf{u}_X \rangle$ and $\langle \mathbf{u}_\beta \rangle$ respectively (and so are equal to $\langle X \rangle$ and $\langle \beta \rangle$). As required, we have reparameterised ϕ_{XY} into a function $\tilde{\phi}_{XY}$ which is a multi-linear function of natural statistic vectors.

3.3 Definition of the Variational Message Passing Algorithm

We have now reached the point where we can specify exactly what form the messages between nodes must take and so define the variational message passing algorithm. The message from a parent node Y to a child node X is just the expectation under Q of the natural statistic vector

$$\mathbf{m}_{Y \rightarrow X} = \langle \mathbf{u}_Y \rangle. \quad (18)$$

The message from a child node X to a parent node Y is

$$\mathbf{m}_{X \rightarrow Y} = \tilde{\phi}_{XY}(\langle \mathbf{u}_X \rangle, \{\mathbf{m}_{i \rightarrow X}\}_{i \in \text{cp}_Y}) \quad (19)$$

which relies on X having received messages previously from all the co-parents. If any node A is observed then the messages are as defined above but with $\langle \mathbf{u}_A \rangle$ replaced by \mathbf{u}_A .

Example If X is Gaussian distributed with conditional $P(X | Y, \beta)$, the messages to its parents Y and β are

$$\mathbf{m}_{X \rightarrow Y} = \begin{bmatrix} \langle \beta \rangle \langle X \rangle \\ -\langle \beta \rangle / 2 \end{bmatrix}, \quad \mathbf{m}_{X \rightarrow \beta} = \begin{bmatrix} -\frac{1}{2} (\langle X^2 \rangle - 2 \langle X \rangle \langle Y \rangle + \langle Y^2 \rangle) \\ \frac{1}{2} \end{bmatrix}$$

and the message from X to any child node is $\begin{bmatrix} \langle X \rangle \\ \langle X^2 \rangle \end{bmatrix}$.

When a node Y has received messages from all parents and children, we can find its updated posterior distribution Q_Y^* by finding its updated natural parameter vector ϕ_Y^* . This vector ϕ_Y^* is computed from all the messages received at a node using

$$\phi_Y^* = \tilde{\phi}_Y(\{\mathbf{m}_{i \rightarrow Y}\}_{i \in \text{pa}_Y}) + \sum_{j \in \text{ch}_Y} \mathbf{m}_{j \rightarrow Y}, \quad (20)$$

which follows from (17). The new expectation of the natural statistic vector $\langle \mathbf{u}_Y \rangle_{Q_Y^*}$ can then be found, as it is a deterministic function of ϕ_Y^* .

The variational message passing algorithm uses these messages to optimise the variational distribution iteratively, as described in Algorithm 1 below. This algorithm requires that the lower bound $\mathcal{L}(Q)$ be evaluated, which will be discussed in Section 3.6.

Algorithm 1 The variational message passing algorithm

1. Initialise each factor distribution Q_j by initialising the corresponding moment vector $\langle \mathbf{u}_j(X_j) \rangle$.
 2. For each node X_j in turn,
 - Retrieve messages from all parent and child nodes, as defined in (18) and (19). This will require child nodes to retrieve messages from the co-parents of X_j .
 - Compute updated natural parameter vector ϕ_j^* using (20).
 - Compute updated moment vector $\langle \mathbf{u}_j(X_j) \rangle$ given the new setting of the parameter vector.
 3. Calculate the new value of the lower bound $\mathcal{L}(Q)$ (if required).
 4. If the increase in the bound is negligible or a specified number of iterations has been reached, stop. Otherwise repeat from step 2.
-

3.4 Example: the Univariate Gaussian Model

To illustrate how variational message passing works, let us apply it to a model which represents a set of observed one-dimensional data $\{x_n\}_{n=1}^N$ with a univariate Gaussian distribution of mean μ and precision γ ,

$$P(\mathbf{x} | \mathcal{H}) = \prod_{n=1}^N \mathcal{N}(x_n | \mu, \gamma^{-1}).$$

We wish to infer the posterior distribution over the parameters μ and γ . In this simple model the exact solution is tractable, which will allow us to compare the approximate posterior with the true posterior. Of course, for any practical application of VMP, the exact posterior would not be tractable otherwise we would not be using approximate inference methods.

In this model, the conditional distribution of each data point x_n is a univariate Gaussian, which is in the exponential family and so its logarithm can be expressed in standard form as

$$\ln P(x_n | \mu, \gamma^{-1}) = \begin{bmatrix} \gamma\mu \\ -\gamma/2 \end{bmatrix}^T \begin{bmatrix} x_n \\ x_n^2 \end{bmatrix} + \frac{1}{2}(\ln \gamma - \gamma\mu^2 - \ln 2\pi)$$

and so $\mathbf{u}_x(x_n) = [x_n, x_n^2]^T$. This conditional can also be written so as to separate out the dependencies on μ and γ

$$\ln P(x_n | \mu, \gamma^{-1}) = \begin{bmatrix} \gamma x_n \\ -\gamma/2 \end{bmatrix}^T \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + \frac{1}{2}(\ln \gamma - \gamma x_n^2 - \ln 2\pi) \tag{21}$$

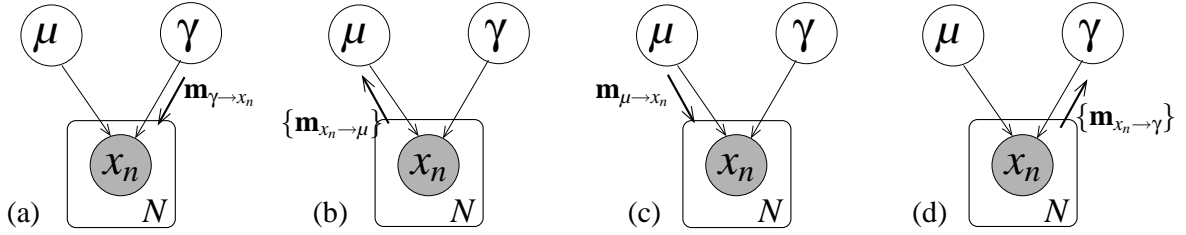


Figure 3: (a)-(d) Message passing procedure for variational inference in a univariate Gaussian model. The box around the x_i node denotes a *plate*, which indicates that the contained node and its connected edges are duplicated N times. The braces around the messages leaving the plate indicate that a set of N distinct messages are being sent.

$$= \begin{bmatrix} -\frac{1}{2}(x_n - \mu)^2 \\ \frac{1}{2} \end{bmatrix}^T \begin{bmatrix} \gamma \\ \ln \gamma \end{bmatrix} - \ln 2\pi \quad (22)$$

which shows that, for conjugacy, $\mathbf{u}_\mu(\mu)$ must be $[\mu, \mu^2]^T$ and $\mathbf{u}_\gamma(\gamma)$ must be $[\gamma, \ln \gamma]^T$ or linear transforms of these.⁴ If we use a separate conjugate prior for each parameter then μ must have a Gaussian prior and γ a gamma prior since these are the exponential family distributions with these natural statistic vectors. Alternatively, we could have chosen a normal-gamma prior over both parameters which leads to a slightly more complicated message passing procedure. We define the parameter priors to have hyper-parameters m, β, a and b , so that

$$\begin{aligned} \ln P(\mu | m, \beta) &= \begin{bmatrix} \beta m \\ -\beta/2 \end{bmatrix}^T \begin{bmatrix} \mu \\ \mu^2 \end{bmatrix} + \frac{1}{2}(\ln \beta - \beta m^2 - \ln 2\pi) \\ \ln P(\gamma | a, b) &= \begin{bmatrix} -b \\ a-1 \end{bmatrix}^T \begin{bmatrix} \gamma \\ \ln \gamma \end{bmatrix} + a \ln b - \ln \Gamma(a). \end{aligned}$$

3.4.1 VARIATIONAL MESSAGE PASSING IN THE UNIVARIATE GAUSSIAN MODEL

We can now apply variational message passing to infer the distributions over μ and γ variationally. The variational distribution is fully factorised and takes the form

$$Q(\mu, \gamma) = Q_\mu(\mu)Q_\gamma(\gamma).$$

We start by initialising $Q_\mu(\mu)$ and $Q_\gamma(\gamma)$ and find initial values of $\langle \mathbf{u}_\mu(\mu) \rangle$ and $\langle \mathbf{u}_\gamma(\gamma) \rangle$. Let us choose to update $Q_\mu(\mu)$ first, in which case variational message passing will proceed as follows (illustrated in Figure 3a-d).

- (a) As we wish to update $Q_\mu(\mu)$, we must first ensure that messages have been sent to the children of μ by any co-parents. Thus, messages $\mathbf{m}_{\gamma \rightarrow x_n}$ are sent from γ to each of the observed nodes x_n . These messages are the same, and are just equal to $\langle \mathbf{u}_\gamma(\gamma) \rangle = [\langle \gamma \rangle, \langle \ln \gamma \rangle]^T$, where the expectation are with respect to the initial setting of Q_γ .

4. To prevent the need for linear transformation of messages, a normalised form of natural statistic vectors will always be used, for example $[\mu, \mu^2]^T$ or $[\gamma, \ln \gamma]^T$.

- (b) Each x_n node has now received messages from all co-parents of μ and so can send a message to μ which is the expectation of the natural parameter vector in (21),

$$\mathbf{m}_{x_n \rightarrow \mu} = \begin{bmatrix} \langle \gamma \rangle x_n \\ -\langle \gamma \rangle / 2 \end{bmatrix}.$$

- (c) Node μ has now received its full complement of incoming messages and can update its natural parameter vector,

$$\Phi_{\mu}^* = \begin{bmatrix} \beta m \\ -\beta / 2 \end{bmatrix} + \sum_{n=1}^N \mathbf{m}_{x_n \rightarrow \mu}.$$

The new expectation $\langle \mathbf{u}_{\mu}(\mu) \rangle$ can then be computed under the updated distribution Q_{μ}^* and sent to each x_n as the message $\mathbf{m}_{\mu \rightarrow x_n} = [\langle \mu \rangle, \langle \mu^2 \rangle]^T$.

- (d) Finally, each x_n node sends a message back to γ which is

$$\mathbf{m}_{x_n \rightarrow \gamma} = \begin{bmatrix} -\frac{1}{2}(x_n^2 - 2x_n \langle \mu \rangle + \langle \mu^2 \rangle) \\ \frac{1}{2} \end{bmatrix}$$

and γ can update its variational posterior

$$\Phi_{\gamma}^* = \begin{bmatrix} -b \\ a - 1 \end{bmatrix} + \sum_{n=1}^N \mathbf{m}_{x_n \rightarrow \gamma}.$$

As the expectation of $\mathbf{u}_{\gamma}(\gamma)$ has changed, we can now go back to step (a) and send an updated message to each x_n node and so on. Hence, in variational message passing, the message passing procedure is repeated again and again until convergence (unlike in belief propagation on a junction tree where the exact posterior is available after a message passing is performed once). Each round of message passing is equivalent to one iteration of the update equations in standard variational inference.

Figure 4 gives an indication of the accuracy of the variational approximation in this model, showing plots of both the true and variational posterior distributions for a toy example. The difference in shape between the two distributions is due to the requirement that Q be factorised. Because $\text{KL}(Q||P)$ has been minimised, the optimal Q is the factorised distribution which lies slightly *inside* P .

3.5 Initialisation and Message Passing Schedule

The variational message passing algorithm is guaranteed to converge to a local minimum of the KL divergence. As with many approximate inference algorithms, including Expectation-Maximisation and Expectation Propagation, it is important to have a good initialisation to ensure that the local minimum that is found is sufficiently close to the global minimum. What makes a good initialisation will depend on the model. In some cases, initialising each factor to a broad distribution will suffice, whilst in others it may be necessary to use a heuristic, such as using K-means to initialise a mixture model.

The variational distribution in the example of Section 3.4 contained only two factors and so messages were passed back-and-forth so as to update these alternately. In fact, unlike belief propagation,

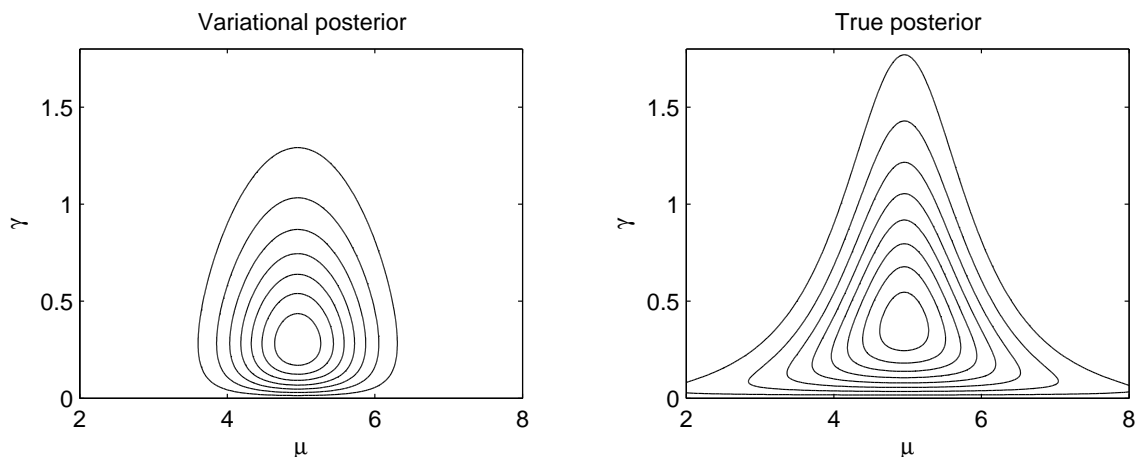


Figure 4: Contour plots of the variational and true posterior over the mean μ and precision γ of a Gaussian distribution, given four samples from $\mathcal{N}(x|5, 1)$. The parameter priors are $P(\mu) = \mathcal{N}(0, 1000)$ and $P(\gamma) = \text{Gamma}(0.001, 0.001)$.

messages in VMP can be passed according to a very flexible schedule. At any point, any factor can be selected and it can be updated locally using only messages from its neighbours and co-parents. There is no requirement that factors be updated in any particular order. However, changing the update order can change which stationary point the algorithm converges to, even if the initialisation is unchanged.

Another constraint on belief propagation is that it is only exact for graphs which are trees and suffers from double-counting if loops are included. VMP does not have this restriction and can be applied to graphs of general form.

3.6 Calculation of the Lower Bound $\mathcal{L}(Q)$

The variational message passing algorithm makes use of the lower bound $\mathcal{L}(Q)$ as a diagnostic of convergence. Evaluating the lower bound is also useful for performing model selection, or model averaging, because it provides an estimate of the log evidence for the model.

The lower bound can also play a useful role in helping to check the correctness both of the analytical derivation of the update equations and of their software implementation, simply by evaluating the bound after updating each factor in the variational posterior distribution and checking that the value of the bound does not decrease. This can be taken a stage further (Bishop and Svensén, 2003) by using numerical differentiation applied to the lower bound. After each update, the gradient of the bound is evaluated in the subspace corresponding to the parameters of the updated factor, to check that it is zero (within numerical tolerances). This requires that the differentiation take account of any constraints on the parameters (for instance that they be positive or that they sum to one). These checks, of course, provide necessary but not sufficient conditions for correctness. Also, they add computational cost so would typically only be employed whilst debugging the implementation.

In previous applications of variational inference, however, the evaluation of the lower bound has typically been done using separate code from that used to implement the update equations.

Although the correctness tests discussed above also provide a check on the mutual consistency of the two bodies of code, it would clearly be more elegant if their evaluation could be unified.

This is achieved naturally in the variational message passing framework by providing a way to calculate the bound automatically, as will now be described. To recap, the lower bound on the log evidence is defined to be

$$\mathcal{L}(Q) = \langle \ln P(\mathbf{H}, \mathbf{V}) \rangle - \langle \ln Q(\mathbf{H}) \rangle,$$

where the expectations are with respect to Q . In a Bayesian network, with a factorised Q distribution, the bound becomes

$$\begin{aligned} \mathcal{L}(Q) &= \sum_i \langle \ln P(X_i | \text{pa}_i) \rangle - \sum_{i \in \mathbf{H}} \langle \ln Q_i(H_i) \rangle \\ &\stackrel{\text{def}}{=} \sum_i \mathcal{L}_i \end{aligned}$$

where it has been decomposed into contributions from the individual nodes $\{\mathcal{L}_i\}$. For a particular latent variable node H_j , the contribution is

$$\mathcal{L}_j = \langle \ln P(H_j | \text{pa}_j) \rangle - \langle \ln Q_j(H_j) \rangle.$$

Given that the model is conjugate-exponential, we can substitute in the standard form for the exponential family

$$\begin{aligned} \mathcal{L}_j &= \langle \phi_j(\text{pa}_j)^T \rangle \langle \mathbf{u}_j(H_j) \rangle + \langle f_j(H_j) \rangle + \langle g_j(\text{pa}_j) \rangle \\ &\quad - \left[\phi_j^{*T} \langle \mathbf{u}_j(H_j) \rangle + \langle f_j(H_j) \rangle + \tilde{g}_j(\phi_j^*) \right], \end{aligned}$$

where the function \tilde{g}_j is a reparameterisation of g_j so as to make it a function of the natural parameter vector rather than the parent variables. This expression simplifies to

$$\mathcal{L}_j = (\langle \phi_j(\text{pa}_j) \rangle - \phi_j^*)^T \langle \mathbf{u}_j(H_j) \rangle + \langle g_j(\text{pa}_j) \rangle - \tilde{g}_j(\phi_j^*). \quad (23)$$

Three of these terms are already calculated during the variational message passing algorithm: $\langle \phi_j(\text{pa}_j) \rangle$ and ϕ_j^* when finding the posterior distribution over H_j in (20), and $\langle \mathbf{u}_j(H_j) \rangle$ when calculating outgoing messages from H_j . Thus, considerable saving in computation are made compared to when the bound is calculated separately.

Each observed variable V_k also makes a contribution to the bound

$$\begin{aligned} \mathcal{L}_k &= \langle \ln P(V_k | \text{pa}_k) \rangle \\ &= \langle \phi_k(\text{pa}_k) \rangle^T \mathbf{u}_k(V_k) + f_k(V_k) + \tilde{g}_k(\langle \phi_k(\text{pa}_k) \rangle). \end{aligned}$$

Again, computation can be saved by computing $\mathbf{u}_k(V_k)$ during the initialisation of the message passing algorithm.

Example 1 *Calculation of the Bound for the Univariate Gaussian Model*

In the univariate Gaussian model, the bound contribution from each observed node x_n is

$$\mathcal{L}_{x_n} = \begin{bmatrix} \langle \gamma \rangle \langle \mu \rangle \\ -\langle \gamma \rangle / 2 \end{bmatrix}^T \begin{bmatrix} x_n \\ x_n^2 \end{bmatrix} + \frac{1}{2} (\langle \ln \gamma \rangle - \langle \gamma \rangle \langle \mu^2 \rangle - \ln 2\pi)$$

and the contributions from the parameter nodes μ and γ are

$$\begin{aligned} \mathcal{L}_\mu &= \begin{bmatrix} \beta m - \beta' m' \\ -\beta/2 + \beta'/2 \end{bmatrix}^T \begin{bmatrix} \langle \mu \rangle \\ \langle \mu^2 \rangle \end{bmatrix} + \frac{1}{2} (\ln \beta - \beta m^2 - \ln \beta' + \beta' m'^2) \\ \mathcal{L}_\gamma &= \begin{bmatrix} -b + b' \\ a - a' \end{bmatrix}^T \begin{bmatrix} \langle \gamma \rangle \\ \langle \ln \gamma \rangle \end{bmatrix} + a \ln b - \ln \Gamma(a) - a' \ln b' + \ln \Gamma(a'). \end{aligned}$$

The bound for this univariate Gaussian model is given by the sum of the contributions from the μ and γ nodes and all x_n nodes.

4. Allowable Models

The variational message passing algorithm can be applied to a wide class of models, which will be characterised in this section.

4.1 Conjugacy Constraints

The main constraint on the model is that each parent–child edge must satisfy the constraint of conjugacy. Conjugacy allows a Gaussian variable to have a Gaussian parent for its mean and we can extend this hierarchy to any number of levels. Each Gaussian node has a gamma parent as the distribution over its precision. Furthermore, each gamma distributed variable can have a gamma distributed scale parameter b , and again this hierarchy can be extended to multiple levels.

A discrete variable can have multiple discrete parents with a Dirichlet prior over the entries in the conditional probability table. This allows for an arbitrary graph of discrete variables. A variable with an Exponential or Poisson distribution can have a gamma prior over its scale or mean respectively, although, as these distributions do not lead to hierarchies, they may be of limited interest.

These constraints are listed in Table 1. This table can be encoded in implementations of the variational message passing algorithm and used during initialisation to check the conjugacy of the supplied model.

4.1.1 TRUNCATED DISTRIBUTIONS

The conjugacy constraint does not put any restrictions on the $f_X(X)$ term in the exponential family distribution. If we choose f_X to be a step function

$$f_X(X) = \begin{cases} 0 & : X \geq 0 \\ -\infty & : X < 0 \end{cases}$$

then we end up with a rectified distribution, so that $P(X | \theta) = 0$ for $X < 0$. The choice of such a truncated distribution will change the form of messages to parent nodes (as the g_X normalisation function will also be different) but will not change the form of messages that are passed to child nodes. However, truncation will affect how the moments of the distribution are calculated from the updated parameters, which will lead to different values of child messages. For example, the moments of a rectified Gaussian distribution are expressed in terms of the standard ‘erf’ function. Similarly, we can consider doubly truncated distributions which are non-zero only over some finite interval, as long as the calculation of the moments and parent messages remains tractable. One

Distribution	1 st parent	Conjugate dist.	2 nd parent	Conjugate dist.
Gaussian	mean μ	Gaussian	precision γ	gamma
gamma	shape a	None	scale b	gamma
discrete	probabilities \mathbf{p}	Dirichlet	parents $\{x_i\}$	discrete
Dirichlet	pseudo-counts \mathbf{a}	None		
Exponential	scale a	gamma		
Poisson	mean $\hat{\lambda}$	gamma		

Table 1: Distributions for each parameter of a number of exponential family distributions if the model is to satisfy conjugacy constraints. Conjugacy also holds if the distributions are replaced by their multivariate counterparts e.g. the distribution conjugate to the precision matrix of a multivariate Gaussian is a Wishart distribution. Where “None” is specified, no standard distribution satisfies conjugacy.

potential problem with the use of a truncated distribution is that no standard distributions may exist which are conjugate for each distribution parameter.

4.2 Deterministic Functions

We can considerably enlarge the class of tractable models if variables are allowed to be defined as deterministic functions of the states of their parent variables. This is achieved by adding deterministic nodes into the graph, as have been used to similar effect in the BUGS software (see Section 5).

Consider a deterministic node X which has stochastic parents $\mathbf{Y} = \{Y_1, \dots, Y_M\}$ and which has a stochastic child node Z . The state of X is given by a deterministic function f of the state of its parents, so that $X = f(\mathbf{Y})$. If X were stochastic, the conjugacy constraint with Z would require that $P(X|\mathbf{Y})$ must have the same functional form, with respect to X , as $P(Z|X)$. This in turn would dictate the form of the natural statistic vector \mathbf{u}_X of X , whose expectation $\langle \mathbf{u}_X(X) \rangle_Q$ would be the message from X to Z .

Returning to the case where X is deterministic, it is still necessary to provide a message to Z of the form $\langle \mathbf{u}_X(X) \rangle_Q$ where the function \mathbf{u}_X is dictated by the conjugacy constraint. This message can be evaluated only if it can be expressed as a function of the messages from the parent variables, which are the expectations of their natural statistics functions $\{\langle \mathbf{u}_{Y_i}(Y_i) \rangle_Q\}$. In other words, there must exist a vector function ψ_X such that

$$\langle \mathbf{u}_X(f(\mathbf{Y})) \rangle_Q = \psi_X(\langle \mathbf{u}_{Y_1}(Y_1) \rangle_Q, \dots, \langle \mathbf{u}_{Y_M}(Y_M) \rangle_Q).$$

As was discussed in Section 3.2, this constrains $\mathbf{u}_X(f(\mathbf{Y}))$ to be a multi-linear function of the set of functions $\{\mathbf{u}_{Y_i}(Y_i)\}$.

A deterministic node can be viewed as having a conditional distribution which is a delta function, so that $P(X|\mathbf{Y}) = \delta(X - f(\mathbf{Y}))$. If X is discrete, this is the distribution that assigns probability one to the state $X = f(\mathbf{Y})$ and zero to all other states. If X is continuous, this is the distribution with the property that $\int g(X) \delta(X - f(\mathbf{Y})) dX = g(f(\mathbf{Y}))$. The contribution to the lower bound from a deterministic node is zero.

Example 2 Using a Deterministic Function as the Mean of a Gaussian

Consider a model where a deterministic node X is to be used as the mean of a child Gaussian distribution $\mathcal{N}(Z|X, \beta^{-1})$ and where X equals a function f of Gaussian-distributed variables Y_1, \dots, Y_M . The natural statistic vectors of X (as dictated by conjugacy with Z) and those of Y_1, \dots, Y_M are

$$\mathbf{u}_X(X) = \begin{bmatrix} X \\ X^2 \end{bmatrix}, \quad \mathbf{u}_{Y_i}(Y_i) = \begin{bmatrix} Y_i \\ Y_i^2 \end{bmatrix} \text{ for } i = 1 \dots M$$

The constraint on f is that $\mathbf{u}_X(f)$ must be multi-linear in $\{\mathbf{u}_{Y_i}(Y_i)\}$ and so both f and f^2 must be multi-linear in $\{Y_i\}$ and $\{Y_i^2\}$. Hence, f can be any multi-linear function of Y_1, \dots, Y_M . In other words, the mean of a Gaussian can be the sum of products of other Gaussian-distributed variables.

Example 3 Using a Deterministic Function as the Precision of a Gaussian

As another example, consider a model where X is to be used as the precision of a child Gaussian distribution $\mathcal{N}(Z|\mu, X^{-1})$ and where X is a function f of gamma-distributed variables Y_1, \dots, Y_M . The natural statistic vectors of X and Y_1, \dots, Y_M are

$$\mathbf{u}_X(X) = \begin{bmatrix} X \\ \ln X \end{bmatrix}, \quad \mathbf{u}_{Y_i}(Y_i) = \begin{bmatrix} Y_i \\ \ln Y_i \end{bmatrix} \text{ for } i = 1 \dots M.$$

and so both f and $\ln f$ must be multi-linear in $\{Y_i\}$ and $\{\ln Y_i\}$. This restricts f to be proportional to a product of the variables Y_1, \dots, Y_M as the logarithm of a product can be found in terms of the logarithms of terms in that product. Hence $f = c \prod_i Y_i$ where c is a constant. A function containing a summation, such as $f = \sum_i Y_i$, would not be valid as the logarithm of the sum cannot be expressed as a multi-linear function of Y_i and $\ln Y_i$.

4.2.1 VALIDATING CHAINS OF DETERMINISTIC FUNCTIONS

The validity of a deterministic function for a node X is dependent on the form of the stochastic nodes it is connected to, as these dictate the functions \mathbf{u}_X and $\{\mathbf{u}_{Y_i}(Y_i)\}$. For example, if the function was a summation $f = \sum_i Y_i$, it would be valid for the first of the above examples but not for the second. In addition, it is possible for deterministic functions to be chained together to form more complicated expressions. For example, the expression $X = Y_1 + Y_2 Y_3$ can be achieved by having a deterministic product node A with parents Y_2 and Y_3 and a deterministic sum node X with parents Y_1 and A . In this case, the form of the function \mathbf{u}_A is not determined directly by its immediate neighbours, but instead is constrained by the requirement of consistency for the connected deterministic subgraph.

In a software implementation of variational message passing, the validity of a particular deterministic structure can most easily be checked by requiring that the function \mathbf{u}_{X_i} be specified explicitly for each deterministic node X_i , thereby allowing the existing mechanism for checking conjugacy to be applied uniformly across both stochastic and deterministic nodes.

4.2.2 DETERMINISTIC NODE MESSAGES

To examine message passing for deterministic nodes, we must consider the general case where the deterministic node X has multiple children $\{Z_j\}$. The message from the node X to any child Z_j is simply

$$\begin{aligned} \mathbf{m}_{X \rightarrow Z_j} &= \langle \mathbf{u}_X(f(\mathbf{Y})) \rangle_Q \\ &= \Psi_X(\mathbf{m}_{Y_1 \rightarrow X}, \dots, \mathbf{m}_{Y_M \rightarrow X}). \end{aligned}$$

For a particular parent Y_k , the function $\mathbf{u}_X(f(\mathbf{Y}))$ is linear with respect to $\mathbf{u}_{Y_k}(Y_k)$ and so it can be written as

$$\mathbf{u}_X(f(\mathbf{Y})) = \Psi_{X,Y_k}(\{\mathbf{u}_{Y_i}(Y_i)\}_{i \neq k}) \cdot \mathbf{u}_{Y_k}(Y_k) + \lambda(\{\mathbf{u}_{Y_i}(Y_i)\}_{i \neq k})$$

where Ψ_{X,Y_k} is a matrix function of the natural statistics vectors of the co-parents of Y_k . The message from a deterministic node to a parent Y_k is then

$$\mathbf{m}_{X \rightarrow Y_k} = \left[\sum_j \mathbf{m}_{Z_j \rightarrow X} \right] \Psi_{X,Y_k}(\{\mathbf{m}_{Y_i \rightarrow X}\}_{i \neq k})$$

which relies on having received messages from all the child nodes and from all the co-parents. The sum of child messages can be computed and stored locally at the node and used to evaluate all child-to-parent messages. In this sense, it can be viewed as the natural parameter vector of a distribution which acts as a kind of pseudo-posterior over the value of X .

4.3 Mixture Distributions

So far, only distributions from the exponential family have been considered. Often it is desirable to use richer distributions that better capture the structure of the system that generated the data. Mixture distributions, such as mixtures of Gaussians, provide one common way of creating richer probability densities. A mixture distribution over a variable X is a weighted sum of a number of component distributions

$$P(X | \{\pi_k\}, \{\theta_k\}) = \sum_{k=1}^K \pi_k P_k(X | \theta_k)$$

where each P_k is a component distribution with parameters θ_k and a corresponding mixing coefficient π_k indicating the weight of the distribution in the weighted sum. The K mixing coefficients must be non-negative and sum to one.

A mixture distribution is not in the exponential family and therefore cannot be used directly as a conditional distribution within a conjugate-exponential model. Instead, we can introduce an additional discrete latent variable λ which indicates from which component distribution each data point was drawn, and write the distribution as

$$P(X | \lambda, \{\theta_k\}) = \prod_{k=1}^K P_k(X | \theta_k)^{\delta_{\lambda k}}.$$

Conditioned on this new variable, the distribution is now in the exponential family provided that all of the component distributions are also in the exponential family. In this case, the log conditional probability of X given all the parents (including λ) can be written as

$$\ln P(X | \lambda, \{\theta_k\}) = \sum_k \delta(\lambda, k) [\phi_k(\theta_k)^T \mathbf{u}_k(X) + f_k(X) + g_k(\theta_k)].$$

If X has a child Z , then conjugacy will require that all the component distributions have the same natural statistic vector, which we can then call \mathbf{u}_X so: $\mathbf{u}_1(X) = \mathbf{u}_2(X) = \dots = \mathbf{u}_K(X) \stackrel{\text{def}}{=} \mathbf{u}_X(X)$. In addition, we may choose to specify, as part of the model, that all these distributions have exactly

the same form (that is, $f_1 = f_2 = \dots = f_K \stackrel{\text{def}}{=} f_X$), although this is not required by conjugacy. In this case, where all the distributions are the same, the log conditional becomes

$$\begin{aligned} \ln P(X | \lambda, \{\theta_k\}) &= \left[\sum_k \delta(\lambda, k) \phi_k(\theta_k) \right]^T \mathbf{u}_X(X) + f_X(X) \\ &\quad + \sum_k \delta(\lambda, k) g_k(\theta_k) \\ &= \phi_X(\lambda, \{\theta_k\})^T \mathbf{u}_X(X) + f_X(X) + \tilde{g}_X(\phi_X(\lambda, \{\theta_k\})) \end{aligned}$$

where we have defined $\phi_X = \sum_k \delta(\lambda, k) \phi_k(\theta_k)$ to be the natural parameter vector of this mixture distribution and the function \tilde{g}_X is a reparameterisation of g_X to make it a function of ϕ_X (as in Section 3.6). The conditional is therefore in the same exponential family form as each of the components.

We can now apply variational message passing. The message from the node X to any child is $\langle \mathbf{u}_X(X) \rangle$ as calculated from the mixture parameter vector $\phi_X(\lambda, \{\theta_k\})$. Similarly, the message from X to a parent θ_k is the message that would be sent by the corresponding component if it were not in a mixture, scaled by the variational posterior over the indicator variable $Q(\lambda = k)$. Finally, the message from X to λ is the vector of size K whose k th element is $\langle \ln P_k(X | \theta_k) \rangle$.

4.4 Multivariate Distributions

Until now, only scalar variables have been considered. It is also possible to handle vector variables in this framework (or to handle scalar variables which have been grouped into a vector to capture posterior dependencies between the variables). In each case, a multivariate conditional distribution is defined in the overall joint distribution P and the corresponding factor in the variational posterior Q will also be multivariate, rather than factorised with respect to the elements of the vector. To understand how multivariate distributions are handled, consider the d -dimensional Gaussian distribution with mean μ and precision matrix⁵ Λ :

$$P(\mathbf{x} | \mu, \Lambda^{-1}) = \sqrt{\frac{|\Lambda|}{(2\pi)^d}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Lambda (\mathbf{x} - \mu)\right).$$

This distribution can be written in exponential family form

$$\ln \mathcal{N}(\mathbf{x} | \mu, \Lambda^{-1}) = \begin{bmatrix} \Lambda \mu \\ -\frac{1}{2} \text{vec}(\Lambda) \end{bmatrix}^T \begin{bmatrix} \mathbf{x} \\ \text{vec}(\mathbf{x}\mathbf{x}^T) \end{bmatrix} + \frac{1}{2}(\ln |\Lambda| - \mu^T \Lambda \mu - d \ln 2\pi)$$

where $\text{vec}(\cdot)$ is a function that re-arranges the elements of a matrix into a column vector in some consistent fashion, such as by concatenating the columns of the matrix. The natural statistic function for a multivariate distribution therefore depends on both the type of the distribution and its dimensionality d . As a result, the conjugacy constraint between a parent node and a child node will also constrain the dimensionality of the corresponding vector-valued variables to be the same. Multivariate conditional distributions can therefore be handled by VMP like any other exponential family distribution, which extends the class of allowed distributions to include multivariate Gaussian and Wishart distributions.

5. The precision matrix of a multivariate Gaussian is the inverse of its covariance matrix.

A group of scalar variables can act as a single parent of a vector-valued node. This is achieved using a deterministic *concatenation* function which simply concatenates a number of scalar values into a vector. In order for this to be a valid function, the scalar distributions must still be conjugate to the multivariate distribution. For example, a set of d univariate Gaussian distributed variables can be concatenated to act as the mean of a d -dimensional multivariate Gaussian distribution.

4.4.1 NORMAL-GAMMA DISTRIBUTION

The mean μ and precision γ parameters of a Gaussian distribution can be grouped together into a single bivariate variable $\mathbf{c} = \{\mu, \gamma\}$. The conjugate distribution for this variable is the normal-gamma distribution, which is written

$$\ln P(\mathbf{c} | m, \lambda, a, b) = \begin{bmatrix} m\lambda \\ -\frac{1}{2}\lambda \\ -b - \frac{1}{2}\lambda m^2 \\ a - \frac{1}{2} \end{bmatrix} \begin{bmatrix} \mu\gamma \\ \mu^2\gamma \\ \gamma \\ \ln \gamma \end{bmatrix} + \frac{1}{2}(\ln \lambda - \ln 2\pi) + a \ln b - \ln \Gamma(a).$$

This distribution therefore lies in the exponential family and can be used within VMP instead of separate Gaussian and gamma distributions. In general, grouping these variables together will improve the approximation and so increase the lower bound. The multivariate form of this distribution, the normal-Wishart distribution, is handled as described above.

4.5 Summary of Allowable Models

In summary, the variational message passing algorithm can handle probabilistic models with the following very general architecture: arbitrary directed acyclic subgraphs of multinomial discrete variables (each having Dirichlet priors) together with arbitrary subgraphs of univariate and multivariate linear Gaussian nodes (having gamma and Wishart priors), with arbitrary mixture nodes providing connections from the discrete to the continuous subgraphs. In addition, deterministic nodes can be included to allow parameters of child distributions to be deterministic functions of parent variables. Finally, any of the continuous distributions can be singly or doubly truncated to restrict the range of allowable values, provided that the appropriate moments under the truncated distribution can be calculated along with any necessary parent messages.

This architecture includes as special cases models such as hidden Markov models, Kalman filters, factor analysers, principal component analysers and independent component analysers, as well as mixtures and hierarchical mixtures of these.

5. VIBES: An Implementation of Variational Message Passing

The variational message passing algorithm has been implemented in a software package called VIBES (Variational Inference in BayEsian networkS), first described by Bishop et al. (2002). Inspired by WinBUGS (a graphical user interface for BUGS by Lunn et al., 2000), VIBES allows for models to be specified graphically, simply by constructing the Bayesian network for the model. This involves drawing the graph for the network (using operations similar to those in a drawing package) and then assigning properties to each node such as its name, the functional form of the conditional distribution, its dimensionality and its parents. As an example, Figure 5 shows the Bayesian network for the univariate Gaussian model along with a screenshot of the same model in

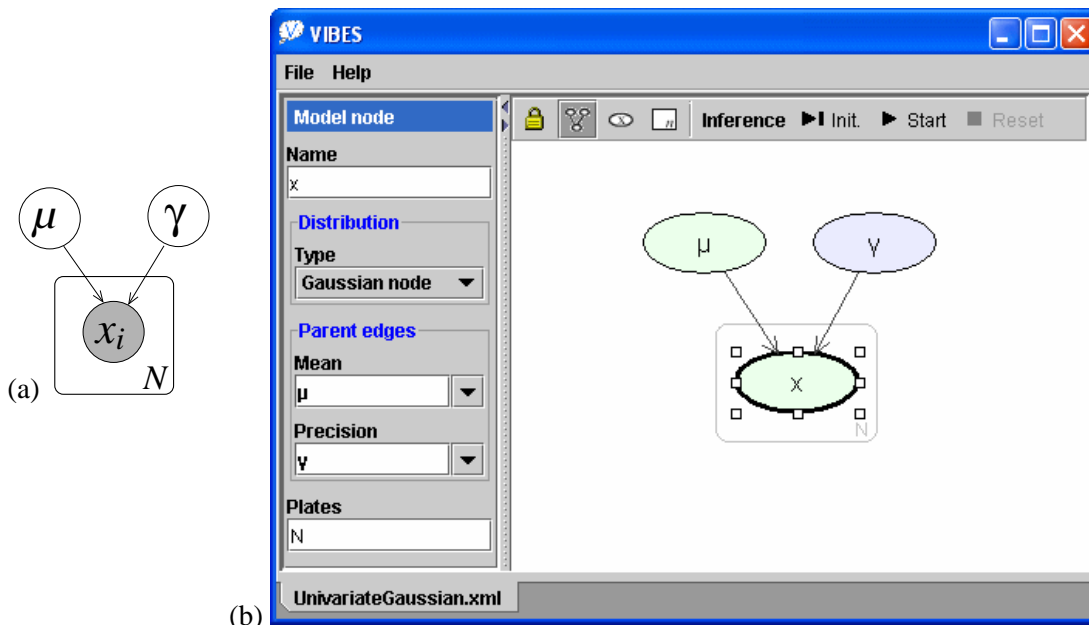


Figure 5: (a) Bayesian network for the univariate Gaussian model. (b) Screenshot of VIBES showing how the same model appears as it is being edited. The node x is selected and the panel to the left shows that it has a Gaussian conditional distribution with mean μ and precision γ . The plate surrounding x shows that it is duplicated N times and the heavy border indicates that it is observed (according to the currently attached data file).

VIBES. Models can also be specified in a text file, which contains XML according to a pre-defined model definition schema. VIBES is written in Java and so can be used on Windows, Linux or any operating system with a Java 1.3 virtual machine.

As in WinBUGS, the convention of making deterministic nodes explicit in the graphical representation has been adopted, as this greatly simplifies the specification and interpretation of the model. VIBES also uses the plate notation of a box surrounding one or more nodes to denote that those nodes are replicated some number of times, specified by the parameter in the bottom right hand corner of the box.

Once the model is specified, data can be attached from a separate data file which contains observed values for some of the nodes, along with sizes for some or all of the plates. The model can then be *initialised* which involves: (i) checking that the model is valid by ensuring that conjugacy and dimensionality constraints are satisfied and that all parameters are specified; (ii) checking that the observed data is of the correct dimensionality; (iii) allocating memory for all moments and messages; (iv) initialisation of the individual distributions Q_i .

Following a successful initialisation, inference can begin immediately. As inference proceeds, the current state of the distribution Q_i for any node can be inspected using a range of diagnostics including tables of values and Hinton diagrams. If desired, the lower bound $\mathcal{L}(Q)$ can be monitored (at the expense of slightly increased computation), in which case the optimisation can be set to

terminate automatically when the change in the bound during one iteration drops below a small value. Alternatively, the optimisation can be stopped after a fixed number of iterations.

The VIBES software can be downloaded from <http://vibes.sourceforge.net>. This software was written by one of the authors (John Winn) whilst a Ph.D. student at the University of Cambridge and is free and open source. Appendix A contains a tutorial for applying VIBES to an example problem involving a Gaussian Mixture model. The VIBES web site also contains an online version of this tutorial.

6. Extensions to Variational Message Passing

In this section, three extensions to the variational message passing algorithm will be described. These extensions are intended to illustrate how the algorithm can be modified to perform alternative inference calculations and to show how the conjugate-exponential constraint can be overcome in certain circumstances.

6.1 Further Variational Approximations: The Logistic Sigmoid Function

As it stands, the VMP algorithm requires that the model be conjugate-exponential. However, it is possible to sidestep the conjugacy requirement by introducing additional variational parameters and approximating non-conjugate conditional distributions by valid conjugate ones. We will now illustrate how this can be achieved using the example of a conditional distribution over a binary variable $x \in 0, 1$ of the form

$$\begin{aligned} P(x|a) &= \sigma(a)^x [1 - \sigma(a)]^{1-x} \\ &= e^{ax} \sigma(-a) \end{aligned}$$

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

is the logistic sigmoid function.

We take the approach of Jaakkola and Jordan (1996) and use a variational bound for the logistic sigmoid function defined as

$$\sigma(a) \geq F(a, \xi) \stackrel{\text{def}}{=} \sigma(\xi) \exp[(a - \xi)/2 + \lambda(\xi)(a^2 - \xi^2)]$$

where $\lambda(\xi) = [1/2 - g(\xi)]/2\xi$ and ξ is a variational parameter. For any given value of a we can make this bound exact by setting $\xi^2 = a^2$. The bound is illustrated in Figure 6 in which the solid curve shows the logistic sigmoid function $\sigma(a)$ and the dashed curve shows the lower bound $F(a, \xi)$ for $\xi = 2$.

We use this result to define a new lower bound $\tilde{\mathcal{L}} \leq \mathcal{L}$ by replacing each expectation of the form $\langle \ln[e^{ax} \sigma(-a)] \rangle$ with its lower bound $\langle \ln[e^{ax} F(-a, \xi)] \rangle$. The effect of this transformation is to replace the logistic sigmoid function with an exponential, therefore restoring conjugacy to the model. Optimisation of each ξ parameter is achieved by maximising this new bound $\tilde{\mathcal{L}}$, leading to the re-estimation equation

$$\xi^2 = \langle a^2 \rangle_Q.$$

It is important to note that, as the quantity $\tilde{\mathcal{L}}$ involves expectations of $\ln F(-a, \xi)$, it is no longer guaranteed to be exact for any value of ξ .

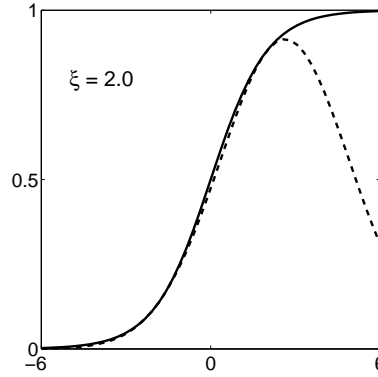


Figure 6: The logistic sigmoid function $\sigma(a)$ and variational bound $F(a, \xi)$.

It follows from (8) that the factor in Q corresponding to $P(x|a)$ is updated using

$$\begin{aligned} \ln Q_x^*(x) &= \langle \ln(e^{ax} F(-a, \xi)) \rangle_{\sim Q_x(x)} + \sum_{k \in \text{ch}_x} \langle \ln P(X_k | \text{pa}_k) \rangle_{\sim Q_x(x)} + \text{const.} \\ &= \langle ax \rangle_{\sim Q_x(x)} + \sum_{k \in \text{ch}_x} \langle b_k x \rangle_{\sim Q_x(x)} + \text{const.} \\ &= a^* x + \text{const.} \end{aligned}$$

where $a^* = \langle a \rangle + \sum_k \langle b_k \rangle$ and the $\{b_k\}$ arise from the child terms which must be in the form ($b_k x + \text{const.}$) due to conjugacy. Therefore, the variational posterior $Q_x(x)$ takes the form

$$Q_x(x) = \sigma(a^*)^x [1 - \sigma(a^*)]^{1-x}.$$

6.1.1 USING THE LOGISTIC APPROXIMATION WITHIN VMP

We will now explain how this additional variational approximation can be used within the VMP framework. The lower bound $\tilde{\mathcal{L}}$ contains terms like $\langle \ln(e^{ax} F(-a, \xi)) \rangle$ which need to be evaluated and so we must be able to evaluate $[\langle a \rangle \langle a^2 \rangle]^T$. The conjugacy constraint on a is therefore that its distribution must have a natural statistic vector $\mathbf{u}_a(a) = [a \ a^2]$. Hence it could, for example, be Gaussian.

For consistency with general discrete distributions, we write the bound on the log conditional $\ln P(x|a)$ as

$$\begin{aligned} \ln P(x|a) &\geq \begin{bmatrix} 0 \\ a \end{bmatrix}^T \begin{bmatrix} \delta(x-0) \\ \delta(x-1) \end{bmatrix} + (-a - \xi)/2 + \lambda(\xi)(a^2 - \xi^2) + \ln \sigma(\xi) \\ &= \begin{bmatrix} \delta(x-1) - \frac{1}{2} \\ \lambda(\xi) \end{bmatrix}^T \begin{bmatrix} a \\ a^2 \end{bmatrix} - \xi/2 - \lambda(\xi)\xi^2 + \ln \sigma(\xi). \end{aligned}$$

The message from node x to node a is therefore

$$\mathbf{m}_{x \rightarrow a} = \begin{bmatrix} \langle \delta(x-1) \rangle - \frac{1}{2} \\ \lambda(\xi) \end{bmatrix}$$

and all other messages are as in standard VMP. The update of variational factors can then be carried out as normal except that each ξ parameter must also be re-estimated during optimisation. This

can be carried out, for example, just before sending a message from x to a . The only remaining modification is to the calculation of the lower bound in (23), where the term $\langle g_j(\text{pa}_j) \rangle$ is replaced by the expectation of its bound,

$$\langle g_j(\text{pa}_j) \rangle \geq (-\langle a \rangle - \xi)/2 + \lambda(\xi)(\langle a^2 \rangle - \xi^2) + \ln \sigma(\xi).$$

This extension to VMP enables discrete nodes to have continuous parents, further enlarging the class of allowable models. In general, the introduction of additional variational parameters enormously extends the class of models to which VMP can be applied, as the constraint that the model distributions must be conjugate no longer applies.

6.2 Finding a Maximum A Posteriori Solution

The advantage of using a variational distribution is that it provides a posterior distribution over latent variables. It is, however, also possible to use VMP to find a Maximum A Posteriori (MAP) solution, in which values of each latent variable are found that maximise the posterior probability. Consider choosing a variational distribution which is a delta function

$$Q^{\text{MAP}}(\mathbf{H}) = \delta(\mathbf{H} - \mathbf{H}^*)$$

where H^* is the MAP solution. From (3), the lower bound is

$$\begin{aligned} \mathcal{L}(Q) &= \langle \ln P(\mathbf{H}, \mathbf{V}) \rangle - \langle \ln Q(\mathbf{H}) \rangle \\ &= \ln P(\mathbf{H}^*, \mathbf{V}) + h_\delta \end{aligned}$$

where h_δ is the differential entropy of the delta function. By considering the differential entropy of a Gaussian in the limit as the variance goes to 0, we can see that $h_\delta = \log a, a \rightarrow 0$. Thus h_δ does not depend on \mathbf{H}^* and so maximising $\mathcal{L}(Q)$ is equivalent to finding the MAP solution. However, since the entropy h_δ tends to $-\infty$, so does $\mathcal{L}(Q)$ and so, whilst it is still trivially a lower bound on the log evidence, it is not an informative one. In other words, knowing the probability density of the posterior at a point is uninformative about the posterior mass.

The variational distribution can be written in factorised form as

$$Q^{\text{MAP}}(\mathbf{H}) = \prod_j Q_j(H_j).$$

with $Q_j(H_j) = \delta(H_j - H_j^*)$. The KL divergence between the approximating distribution and the true posterior is minimised if $\text{KL}(Q_j || Q_j^*)$ is minimised, where Q_j^* is the standard variational solution given by (6). Normally, Q_j is unconstrained so we can simply set it to Q_j^* . However, in this case, Q_j is a delta function and so we have to find the value of H_j^* that minimises $\text{KL}(\delta(H_j - H_j^*) || Q_j^*)$. Unsurprisingly, this is simply the value of H_j^* that maximises $Q_j^*(H_j^*)$.

In the message passing framework, a MAP solution can be obtained for a particular latent variable H_j directly from the updated natural statistic vector ϕ_j^* using

$$(\phi_j^*)^T \frac{d\mathbf{u}_j(H_j)}{dH_j} = 0.$$

For example, if Q_j^* is Gaussian with mean μ then $H_j^* = \mu$ or if Q_j^* is gamma with parameters a, b , then $H_j^* = (a - 1)/b$.

Given that the variational posterior is now a delta function, the expectation of any function $\langle f(H_j) \rangle$ under the variational posterior is just $f(H_j^*)$. Therefore, in any outgoing messages, $\langle \mathbf{u}_j(H_j) \rangle$ is replaced by $\mathbf{u}_j(H_j^*)$. Since all surrounding nodes can process these messages as normal, a MAP solution may be obtained for any chosen subset of variables (such as particular hyper-parameters), whilst a full posterior distribution is retained for all other variables.

6.3 Learning Non-conjugate Priors by Sampling

For some exponential family distribution parameters, there is no standard probability distribution which can act as a conjugate prior. For example, there is no standard distribution which can act as a conjugate prior for the shape parameter a of the gamma distribution. This implies that we cannot learn a posterior distribution over a gamma shape parameter within the basic VMP framework. As discussed above, we can sometimes introduce conjugate approximations by adding variational parameters, but this may not always be possible.

The purpose of the conjugacy constraint is two-fold. First, it means that the posterior distribution of each variable, conditioned on its neighbours, has the same form as the prior distribution. Hence, the updated variational distribution factor for that variable has the same form and inference involves just updating the parameters of that distribution. Second, conjugacy results in variational distributions being in standard exponential family form allowing their moments to be calculated analytically.

If we ignore the conjugacy constraint, we get non-standard posterior distributions and we must resort to using sampling or other methods to determine the moments of these distributions. The disadvantages of using sampling include computational expense, inability to calculate an analytical lower bound and the fact that inference is no longer deterministic for a given initialisation and ordering. The use of sampling methods will now be illustrated by an example showing how to sample from the posterior over the shape parameter of a gamma distribution.

Example 4 Learning a Gamma Shape Parameter

Let us assume that there is a latent variable a which is to be used as the shape parameter of K gamma distributed variables $\{x_1 \dots x_K\}$. We choose a to have a non-conjugate prior of an inverse-gamma distribution:

$$P(a | \alpha, \beta) \propto a^{-\alpha-1} \exp\left(\frac{-\beta}{a}\right).$$

The form of the gamma distribution means that messages sent to the node a are with respect to a natural statistic vector

$$\mathbf{u}_a = \begin{bmatrix} a \\ \ln \Gamma(a) \end{bmatrix}$$

which means that the updated factor distribution Q_a^* has the form

$$\ln Q_a^*(a) = \left[\sum_{i=1}^K \mathbf{m}_{x_i \rightarrow a} \right]^T \begin{bmatrix} a \\ \ln \Gamma(a) \end{bmatrix} + (-\alpha - 1) \ln a - \frac{\beta}{a} + \text{const.}$$

This density is not of standard form, but it can be shown that $Q^*(\ln a)$ is log-concave, so we can generate independent samples from the distribution for $\ln a$ using Adaptive Rejection Sampling from Gilks and Wild (1992). These samples are then transformed to get samples of a from $Q_a^*(a)$, which

is used to estimate the expectation $\langle \mathbf{u}_a(a) \rangle$. This expectation is then sent as the outgoing message to each of the child nodes.

Each factor distribution is normally updated during every iteration and so, in this case, a number of independent samples from Q_a^* would have to be drawn during every iteration. If this proved too computationally expensive, then the distribution need only be updated intermittently.

It is worth noting that, as in this example, BUGS also uses Adaptive Rejection Sampling for sampling when the posterior distribution is log-concave but non-conjugate, whilst also providing techniques for sampling when the posterior is not log-concave. This suggests that non-conjugate parts of a general graphical model could be handled within a BUGS-style framework whilst variational message passing is used for the rest of the model. The resulting hybrid variational/sampling framework would, to a certain extent, capture the advantages of both techniques.

7. Discussion

The variational message passing algorithm allows approximate inference using a factorised variational distribution in any conjugate-exponential model, and in a range of non-conjugate models. As a demonstration of its utility, this algorithm has already been used to solve problems in the domain of machine vision and bioinformatics (see Winn, 2003; Bishop and Winn, 2000). In general, variational message passing dramatically simplifies the construction and testing of new variational models and readily allows a range of alternative models to be tested on a given problem.

The general form of VMP also allows the inclusion of arbitrary nodes in the graphical model provided that each node is able to receive and generate appropriate messages in the required form, whether or not the model remains conjugate-exponential. The extensions to VMP concerning the logistic function and sampling illustrate this flexibility.

One limitation of the current algorithm is that it uses a variational distribution which is factorised across nodes, giving an approximate posterior which is separable with respect to individual (scalar or vector) variables. In general, an improved approximation will be achieved if a posterior distribution is used which retains some dependency structure. Whilst Wiegerinck (2000) has presented a general framework for such structured variational inference, he does not provide a general-purpose algorithm for applying this framework. Winn (2003) and Bishop and Winn (2003) have therefore proposed an extended version of variational message passing which allows for structured variational distributions. VIBES has been extended to implement a limited version of this algorithm that can only be applied to a constrained set of models. However, a complete implementation and evaluation of this extended algorithm has yet to be undertaken.

The VIBES software is free and open source and can be downloaded from the VIBES web site at <http://vibes.sourceforge.net>. The web site also contains a tutorial that provides an introduction to using VIBES.

Acknowledgments

The authors would like to thank David Spiegelhalter for his help with the VIBES project. We would also like to thank Zoubin Ghahramani, David MacKay, Matthew Beal and Michael Jordan for many helpful discussions about variational inference.

This work was carried out whilst John Winn was a Ph.D. student at the University of Cambridge, funded by a Microsoft Research studentship.

Appendix A. VIBES Tutorial

In this appendix, we demonstrate the application of VIBES to an example problem involving a Gaussian Mixture model. We then demonstrate the flexibility of VIBES by changing the model to fit the data better, using the lower bound as an estimate of the log evidence for each model. An online version of this tutorial is available at <http://vibes.sourceforge.net/tutorial>.

The data used in this tutorial is two-dimensional and consists of nine clusters in a three-by-three grid, as illustrated in Figure 7.

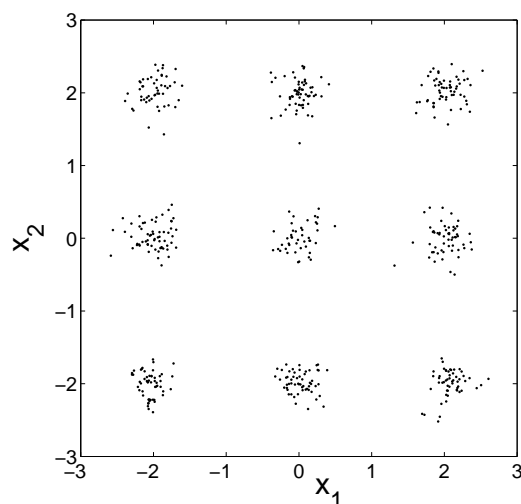


Figure 7: The two-dimensional data set used in the tutorial, which consists of nine clusters in a three-by-three grid.

A.1 Loading Matlab Data into VIBES

The first step is to load the data set into VIBES. This is achieved by creating a node with the name x which corresponds to a matrix x in a Matlab `.mat` file. As the data matrix is two dimensional, the node is placed inside two plates N and d and the data filename (in this case `MixGaussianData2D.mat`) is entered. Selecting `File`→`Load data` loads the data into the node and also sets the size of the N and d plates to 500 and 2 respectively. The node is marked as observed (shown with a bold edge) and the observed data can be inspected by double-clicking the node with the mouse. At this point, the display is as shown in Figure 8.

A.2 Creating and Learning a Gaussian Model

The node x has been marked as Gaussian by default and so the model is invalid as neither the mean nor the precision of the Gaussian have been set (attempting to initialise the model by pressing the `Init.` button will give an error message to this effect). We can specify latent variables for these

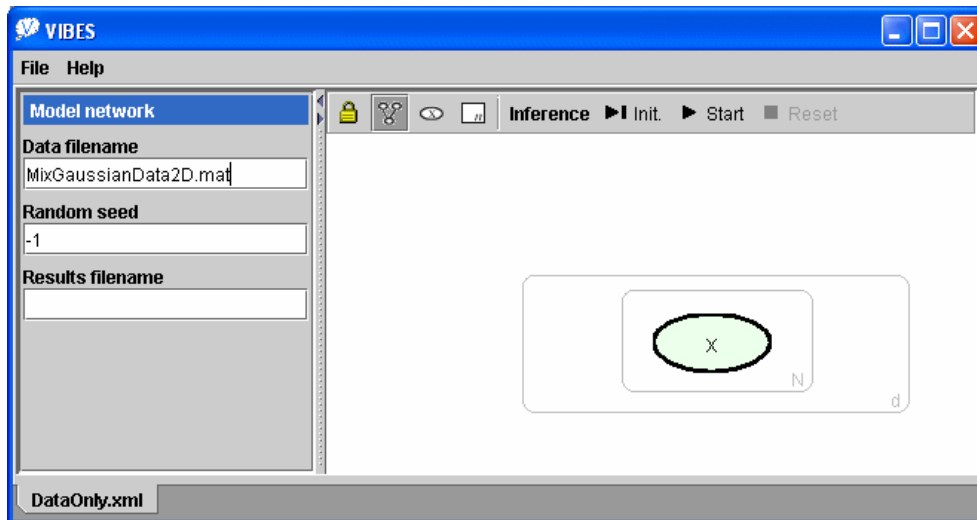


Figure 8: A VIBES model with a single observed node x which has attached data.

parameters by creating a node μ for the mean parameter and a node γ for the precision parameter. These nodes are created within the d plate to give a model which is separable over each data dimension. These are then set as the Mean and Precision properties of x , as shown in Figure 9.

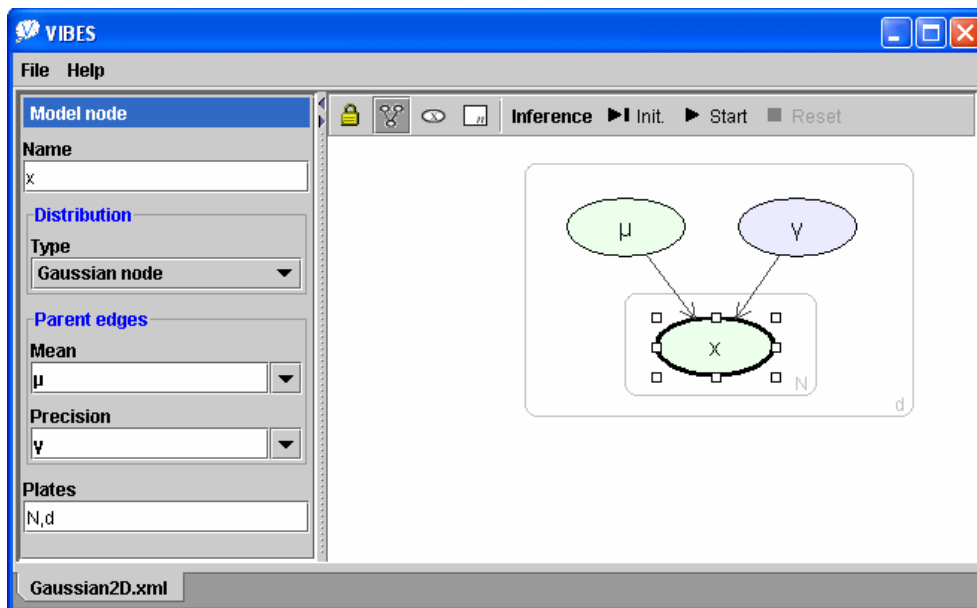


Figure 9: A two-dimensional Gaussian model, showing that the variables μ and γ are being used as the mean and precision parameters of the conditional distribution over x .

The model is still invalid as the parameters of μ and γ are unspecified. In this case, rather than create further latent variables, these parameters will be set to fixed values to give appropriate priors

(for example setting μ to have mean = 0 and precision = 0.3 and γ to have $a = 10$ and $b = 1$). The network now corresponds to a two-dimensional Gaussian model and variational inference can be performed automatically by pressing the `Start` button (which also performs initialisation). For this data set, inference converges after four iterations and gives a bound of -1984 nats. At this point, the expected values of each latent variable under the fully-factorised Q distribution can be displayed or graphed by double-clicking on the corresponding node.

A.3 Extending the Gaussian model to a Gaussian Mixture Model

Our aim is to create a Gaussian mixture model and so we must extend our simple Gaussian model to be a mixture with K Gaussian components. As there will now be K sets of the latent variables μ and γ , these are placed in a new plate, called K , whose size is set to 20. We modify the conditional distribution for the x node to be a mixture of dimension K , with each component being Gaussian. The display is then as shown in Figure 10.

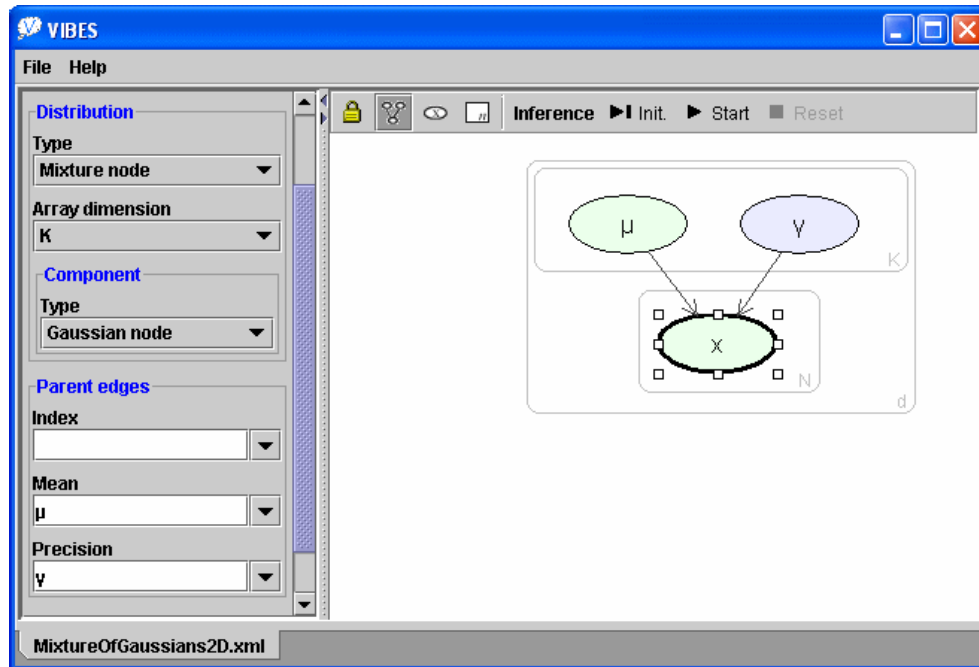


Figure 10: An incomplete model which shows that x is now a mixture of K Gaussians. There are now K sets of parameters and so μ and γ have been placed in a plate K . The model is incomplete as the `Index` parent of x has not been specified.

The model is currently incomplete as making x a mixture requires a new discrete `Index` parent to indicate which component distribution each data point was drawn from. We must therefore create a new node λ , sitting in the N plate, to represent this new discrete latent variable. We also create a node π with a Dirichlet distribution which provides a prior over λ . The completed mixture model is shown in Figure 11.

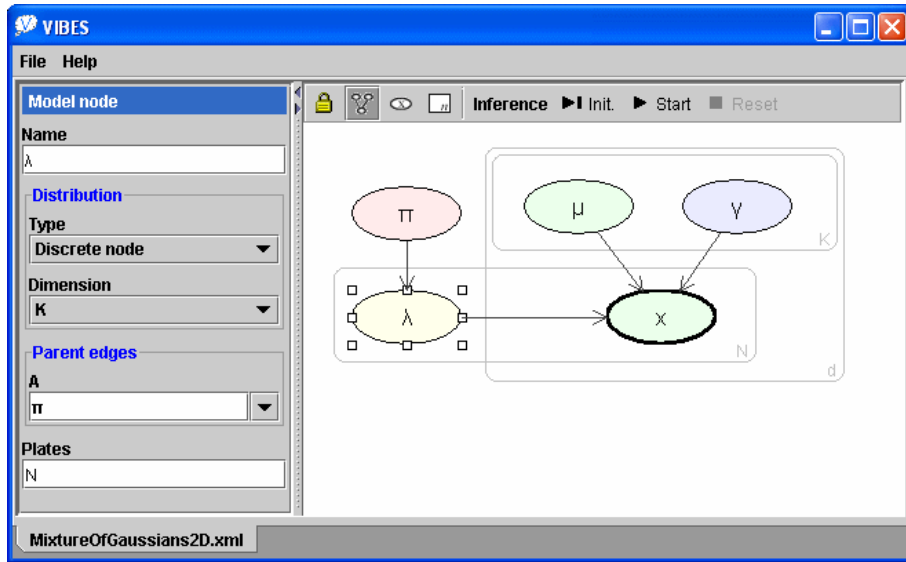


Figure 11: The completed Gaussian mixture model showing the discrete indicator node λ .

A.4 Inference Using the Gaussian Mixture Model

With the model complete, inference can once again proceed automatically by pressing the Start button. A Hinton diagram of the expected value of π can be displayed by double-clicking on the π node, giving the result shown in Figure 12. As can be seen, nine of the twenty components have been retained.

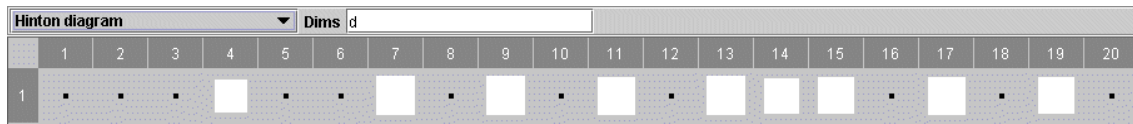


Figure 12: A Hinton diagram showing the expected value of π for each mixture component. The learned mixture consists of only nine components.

The means of the retained components can be inspected by double-clicking on the μ node, giving the Hinton diagram of Figure 13. These learned means correspond to the centres of each of the data clusters.

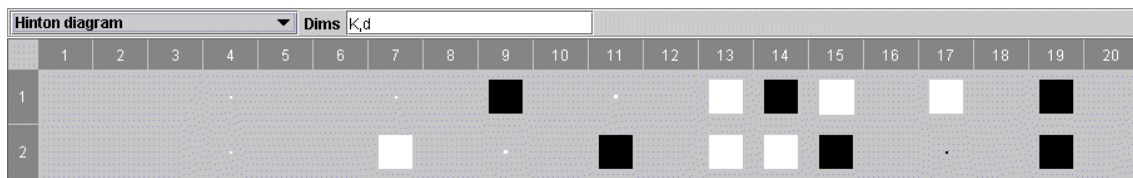


Figure 13: A Hinton diagram whose columns give the expected two-dimensional value of the mean μ for each mixture component. The mean of each of the eleven unused components is just the expected value under the prior which is $(0,0)$. Column 4 corresponds to a retained component whose mean is roughly $(0,0)$.

A graph of the evolution of the bound can be displayed by clicking on the bound value and is shown in Figure 14. The converged lower bound of this new model is -1019 nats, which is significantly higher than that of the single Gaussian model, showing that there is much greater evidence for this model. This is unsurprising since a mixture of 20 Gaussians has significantly more parameters than a single Gaussian and hence can give a much closer fit to the data. Note, however, that the model automatically chooses only to exploit 9 of these components, with the remainder being suppressed (by virtue of their mixing coefficients going to zero). This provides an elegant example of automatic model complexity selection within a Bayesian setting.

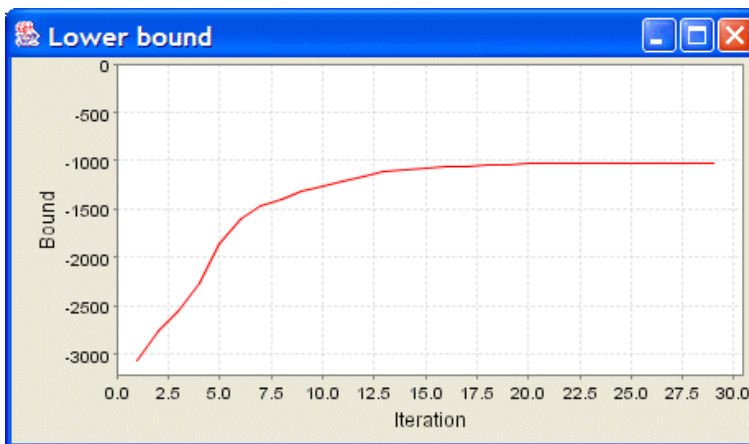


Figure 14: A graph of the evolution of the lower bound during inference.

A.5 Modifying the Mixture Model

The rapidity with which models can be constructed using VIBES allows new models to be quickly developed and compared. For example, we can take our existing mixture of Gaussians model and modify it to try and find a more probable model.

First, we may hypothesise that each of the clusters has similar size and so they may be modelled by a mixture of Gaussian components having a common variance in each dimension. Graphically, this corresponds to shrinking the K plate so that it no longer contains the γ node, as shown in Figure 15a. The converged lower bound for this new model is -937 nats showing that this modified model is better at explaining this data set than the standard mixture of Gaussians model. Note that the increase in model probability does not arise from an improved fit to the data, since this model and the previous one both contain 20 Gaussian components and in both cases 9 of these components contribute to the data fit. Rather, the constrained model having a single variance parameter can achieve almost as good a data fit as the unconstrained model yet with far fewer parameters. Since a Bayesian approach automatically penalises complexity, the simpler (constrained) model has the higher probability as indicated by the higher value for the variational lower bound.

We may further hypothesise that the data set is separable with respect to its two dimensions (i.e. the two dimensions are independent). Graphically this consists of moving all nodes inside the d plate (so we effectively have two copies of a one-dimensional mixture of Gaussians model with common variance). A VIBES screenshot of this further modification is shown in Figure 15b.

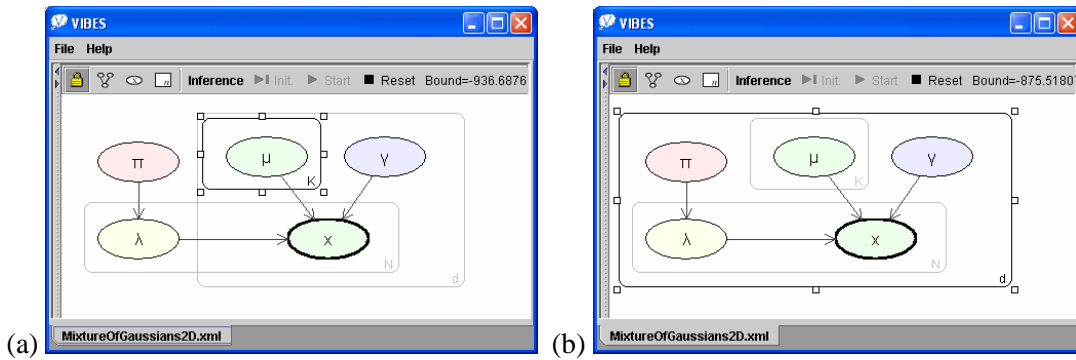


Figure 15: (a) Mixture of Gaussians model with shared precision parameter γ (the γ node is no longer inside the K plate). (b) Model with independent data dimensions, each a univariate Gaussian mixture with common variance.

Performing variational inference on this separable model leads to each one-dimensional mixture having three retained mixture components and gives an improved bound of -876 nats.

We will consider one final model. In this model both the π and the γ nodes are common to both data dimensions, as shown in Figure 16. This change corresponds to the assumption that the mixture coefficients are the same for each of the two mixtures and that the component variances are the same for all components in both mixtures. Inference leads to a final improved bound of -856 nats. Whilst this tutorial has been on a toy data set, the principles of model construction, modification and comparison can be applied just as readily to real data sets.

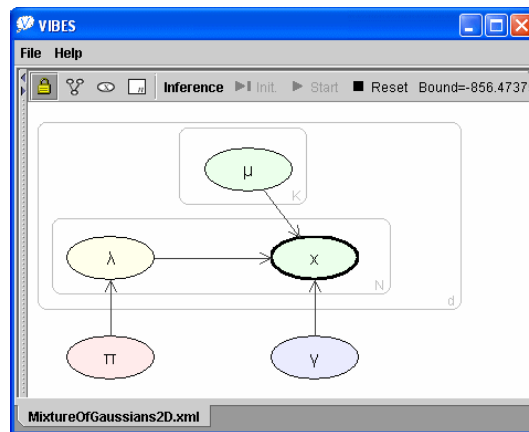


Figure 16: Further modified mixture model where the π and γ nodes are now common to all data dimensions.

References

- H. Attias. A variational Bayesian framework for graphical models. In S. Solla, T. K. Leen, and K-L Muller, editors, *Advances in Neural Information Processing Systems*, volume 12, pages 209–215, Cambridge MA, 2000. MIT Press.
- C. M. Bishop. Variational principal components. In *Proceedings Ninth International Conference on Artificial Neural Networks, ICANN'99*, volume 1, pages 509–514. IEE, 1999.
- C. M. Bishop and M. Svensén. Bayesian Hierarchical Mixtures of Experts. In U. Kjaerulff and C. Meek, editors, *Proceedings Nineteenth Conference on Uncertainty in Artificial Intelligence*, pages 57–64. Morgan Kaufmann, 2003.
- C. M. Bishop and J. M. Winn. Non-linear Bayesian image modelling. In *Proceedings Sixth European Conference on Computer Vision*, volume 1, pages 3–17. Springer-Verlag, 2000.
- C. M. Bishop and J. M. Winn. Structured variational distributions in VIBES. In *Proceedings Artificial Intelligence and Statistics*, Key West, Florida, 2003. Society for Artificial Intelligence and Statistics.
- C. M. Bishop, J. M. Winn, and D. Spiegelhalter. VIBES: A variational inference engine for Bayesian networks. In *Advances in Neural Information Processing Systems*, volume 15, 2002.
- R. G. Cowell, A. P. Dawid, S. L. Lauritzen, and D. J. Spiegelhalter. *Probabilistic Networks and Expert Systems*. Statistics for Engineering and Information Science. Springer-Verlag, 1999.
- Z. Ghahramani and M. J. Beal. Propagation algorithms for variational Bayesian learning. In T. K. Leen, T. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems*, volume 13, Cambridge MA, 2001. MIT Press.
- W. R. Gilks and P. Wild. Adaptive rejection sampling for Gibbs sampling. *Applied Statistics*, 41(2): 337–348, 1992.
- T. Jaakkola and M. Jordan. A variational approach to Bayesian logistic regression problems and their extensions. In *In Proceedings of the 6th international workshop on artificial intelligence and statistics.*, 1996.
- M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An introduction to variational methods for graphical models. In M. I. Jordan, editor, *Learning in Graphical Models*, pages 105–162. Kluwer, 1998.
- S. L. Lauritzen. Propagation of probabilities, means, and variances in mixed graphical association models. *Journal of the American Statistical Association*, 87(420):1098–1108, 1992.
- D. J. Lunn, A. Thomas, N. G. Best, and D. J. Spiegelhalter. WinBUGS – a Bayesian modelling framework: concepts, structure and extensibility. *Statistics and Computing*, 10:321–333, 2000. <http://www.mrc-bsu.cam.ac.uk/bugs/>.
- T. P. Minka. Expectation propagation for approximate Bayesian inference. In *Proceedings of the 17th Annual Conference on Uncertainty in Artificial Intelligence*, pages 362–369. Morgan Kaufmann, 2001.

- R. M. Neal and G. E. Hinton. A new view of the EM algorithm that justifies incremental and other variants. In M. I. Jordan, editor, *Learning in Graphical Models*, pages 355–368. Kluwer, 1998.
- J. Pearl. Fusion, propagation and structuring in belief networks. *Artificial Intelligence*, 29:241–288, 1986.
- A. Thomas, D. J. Spiegelhalter, and W. R. Gilks. BUGS: A program to perform Bayesian inference using Gibbs sampling. In J. M. Bernardo, J. O. Berger, A. P. Dawid, and A. F. M. Smith, editors, *Bayesian Statistics*, Oxford: Clarendon Press, 1992.
- W. Wiegerinck. Variational approximations between mean field theory and the junction tree algorithm. In *Uncertainty in Artificial Intelligence*. Morgan Kauffmann, 2000.
- J. M. Winn. *Variational Message Passing and its Applications*. PhD thesis, University of Cambridge, October 2003.