

Variational inference in (sparse) latent LDS models

bcseke@inf.ed.ac.uk

January 17, 2013

Abstract

Factored expectation constraints based approximate inference for Latent LDS models. Details of sparse models are also considered.

Model

We consider the stationary homogeneous latent Gaussian model

$$\begin{aligned} \mathbf{x}_{t+1} &= \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t + \mathbf{Q}^{-1/2}\mathbf{w}_t, \quad \text{with } \mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \\ p(\mathbf{y}_{t+1}|\mathbf{x}_{t+1}) &= \phi_{t+1,j}(\mathbf{x}_{t+1}; \mathbf{y}_{t+1}) \end{aligned}$$

where the inputs $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_T]$ and the outputs $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_T]$ are observed and known and $\phi_{t+1,j}$ is some, possibly non-Gaussian observation model. We put independent priors on the elements of the parameters \mathbf{A} , \mathbf{B} and \mathbf{Q} and write the joint density conditioned by the inputs as

$$p(\mathbf{Y}, \mathbf{X}, \mathbf{A}, \mathbf{Q}, \mathbf{B}|\mathbf{U}) = \left[\prod_{ij} p_0(a_{ij}) \prod_{ij} p_0(b_{ij}) \prod_{ij} p_0(q_{ij}) \right] \times \prod_t N(\mathbf{x}_{t+1}|\mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t, \mathbf{Q}^{-1}) \prod_j \phi_{t+1,j}(\mathbf{x}_{t+1}; \mathbf{y}_{t+1}).$$

In the following we show how we plan to do inference with a factored expectation propagation algorithm that is expected to get closer to the free form variational inference than the fixed form variational approach. We start from the free form variational approach and then we present the details of the expectation constraints based procedure.

Free form variational inference

Here we choose to approximate the joint density $p(\mathbf{X}, \mathbf{A}, \mathbf{Q}, \mathbf{B}|\mathbf{Y}, \mathbf{U})$ with a factored form $q_x(\mathbf{X})q_A(\mathbf{A})q_Q(\mathbf{Q})q_B(\mathbf{B})$ by using the KL divergence $D[\cdot|\cdot]$, that is we have

$$\text{minimise}_{q_x, q_A, q_Q, q_B} D[q_x(\mathbf{X})q_A(\mathbf{A})q_Q(\mathbf{Q})q_B(\mathbf{B})||p(\mathbf{X}, \mathbf{A}, \mathbf{Q}, \mathbf{B}|\mathbf{Y}, \mathbf{U})].$$

The stationary point of this minimisation problem yield the structured variational mean field updates

$$[q_x(\mathbf{X})]^{new} \propto \exp \left\{ \langle \log p(\mathbf{X}, \mathbf{A}, \mathbf{Q}, \mathbf{B}|\mathbf{Y}, \mathbf{U}) \rangle_{q_A q_Q, q_B} \right\} \quad (1)$$

$$[q_A(\mathbf{A})]^{new} \propto \exp \left\{ \langle \log p(\mathbf{X}, \mathbf{A}, \mathbf{Q}, \mathbf{B}|\mathbf{Y}, \mathbf{U}) \rangle_{q_x q_Q, q_B} \right\} \quad (2)$$

$$[q_Q(\mathbf{Q})]^{new} \propto \exp \left\{ \langle \log p(\mathbf{X}, \mathbf{A}, \mathbf{Q}, \mathbf{B}|\mathbf{Y}, \mathbf{U}) \rangle_{q_x q_A, q_B} \right\} \quad (3)$$

$$[q_B(\mathbf{B})]^{new} \propto \exp \left\{ \langle \log p(\mathbf{X}, \mathbf{A}, \mathbf{Q}, \mathbf{B}|\mathbf{Y}, \mathbf{U}) \rangle_{q_x q_A, q_Q} \right\} \quad (4)$$

From now on we will use and alternative notation for expectations: and expectation $h_i(x_i) \equiv \langle h(x_1, \dots, x_i, \dots, x_n) \rangle_{q_1, \dots, q_{i-1}, q_{i+1}, \dots, q_n}$ is also denoted as $h_i(x_i) = \langle h(x_1, \dots, x_i, \dots, x_n) \rangle_{q_i}$. It

will typically be clear from the context what the densities q_1, \dots, q_n are and which member of this set of densities is omitted.

In the following sections we detail the form of all these densities and models in turn and evaluate what further steps can be taken to approximate them. The crucial term is the transition term that can be written as

$$\log N(\mathbf{x}_{t+1} | \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t, \mathbf{Q}^{-1}) = \frac{1}{2} \log \det(\mathbf{Q}/2\pi) - \frac{1}{2} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t - \mathbf{B}\mathbf{u}_t)^T \mathbf{Q} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t - \mathbf{B}\mathbf{u}_t) \quad (5)$$

As we can see, the above separation of the variables is really justified: due to the form of the interaction there joint modelling of the state space and the parameters is problematic as there is hardly any parametric class that could deal with such high degree of interaction, whereas when considered conditionally, they are all quadratic.

The model for q_x

From (5) we obtain a q_x that can be rewritten as

$$q_x(\mathbf{X}) \propto \prod_t \Psi_{t,t+1}(\mathbf{x}_t, \mathbf{x}_{t+1}) \prod_j \phi_{t+1,j}(\mathbf{x}_{t+1}; \mathbf{y}_{t+1}) \quad (6)$$

with $\Psi_{t,t+1}(\mathbf{x}_{t+1}, \mathbf{x}_t) = \langle \log N(\mathbf{x}_{t+1} | \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t, \mathbf{Q}^{-1}) \rangle_{q_A q_B q_Q}$ is a canonical Gaussian with parameters $\langle \mathbf{h}_{t,t+1}^x \rangle_{q_x}$ and $\langle \mathbf{Q}_{t,t+1}^x \rangle_{q_x}$ where

$$\mathbf{h}_{t,t+1}^x = \begin{bmatrix} -\mathbf{A}^T \mathbf{Q} \mathbf{B} \mathbf{u}_t \\ \mathbf{0} \end{bmatrix}$$

$$\mathbf{Q}_{t,t+1}^x = \begin{bmatrix} \mathbf{A}^T \mathbf{Q} \mathbf{A} & -\mathbf{A}^T \mathbf{Q} \\ -\mathbf{Q} \mathbf{A} & \mathbf{Q} \end{bmatrix}.$$

Therefore, we have a latent Gaussian model and inference in this model can be addressed by well-known approximations.

The model for q_A

Let us use the notation $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$ and let $[\mathbf{A}]_c$ be the column ordered vectorised form of \mathbf{A} , namely $[\mathbf{A}]_c^T = [\mathbf{a}_1^T, \dots, \mathbf{a}_n^T]$. Then it follows that

$$q_A(\mathbf{A}) \propto \exp \left\{ \langle \mathbf{h}_A \rangle_{q_A}^T [\mathbf{A}]_c - \frac{1}{2} [\mathbf{A}]_c^T \langle \mathbf{Q}_A \rangle_{q_A} [\mathbf{A}]_c \right\} \times \prod_{ij} p_0(a_{ij}), \quad (7)$$

where

$$\mathbf{h}_A = \left[\mathbf{Q} \sum_t \mathbf{x}_{t+1} \mathbf{x}_t^T - \mathbf{Q} \mathbf{B} \sum_t \mathbf{u}_t \mathbf{x}_t^T \right]_c,$$

$$\mathbf{Q}_A = \sum_t \mathbf{x}_t \mathbf{x}_t^T \otimes \mathbf{Q}$$

The model for q_B

The model for \mathbf{B} follows a similar structure and is given by

$$q_B(\mathbf{B}) \propto \exp \left\{ \langle \mathbf{h}_B \rangle_{q_B}^T [\mathbf{B}]_c - \frac{1}{2} [\mathbf{B}]_c^T \langle \mathbf{Q}_B \rangle_{q_B} [\mathbf{B}]_c \right\} \times \prod_{ij} p_0(b_{ij}) \quad (8)$$

with

$$\mathbf{h}_B = \left[\mathbf{Q} \sum_t (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t) \mathbf{u}_t^T \right]_c,$$

$$\mathbf{Q}_B = \sum_t \mathbf{u}_t \mathbf{u}_t^T \otimes \mathbf{Q}.$$

The model for q_Q

This model is rather complicated because it involves a log determinant term, but suitable choices for the structure of \mathbf{Q} and the prior $p_0(\mathbf{Q})$ can help to make inference tractable. The distribution q_Q has the form

$$q_Q(\mathbf{Q}) \propto \exp \left\{ \frac{1}{2} \log \det \mathbf{Q} - \frac{1}{2} \text{tr}(\mathbf{Q}^T \mathbf{H}_Q) \right\} \times \prod_{ij} p_0(q_{ij}), \quad (9)$$

where

$$\begin{aligned} \mathbf{H}_q = & \langle \mathbf{x}_{t+1} \mathbf{x}_{t+1} \rangle_{q_x} - 2 \text{tr} \left\{ \langle \mathbf{A} \rangle_{q_A} \langle \mathbf{x}_t \mathbf{x}_{t+1}^T \rangle_{q_x} \right\} + \text{tr} \left\{ \langle \mathbf{A}^T \mathbf{A} \rangle_{q_A} \langle \mathbf{x}_t \mathbf{x}_t^T \rangle_{q_x} \right\} \\ & - 2 \langle \mathbf{B} \rangle_{q_B} \sum_t \mathbf{u}_t (\langle \mathbf{x}_{t+1} \rangle_{q_x} - \langle \mathbf{A} \rangle_{q_A} \langle \mathbf{x}_t \rangle_{q_x})^T + \langle \mathbf{B}^T \mathbf{B} \rangle_{q_B} \sum_t \mathbf{u}_t \mathbf{u}_t^T. \end{aligned}$$

Inference with factored expectation propagation

With the exception of q_Q , all model above are latent Gaussian models where exact inference can be intractable due to the priors or the likelihoods. For this reason we will introduce a factored EP algorithm which basically simplifies to doing EP in these models. In the following we introduce a family of marginals, expectation constraints and the corresponding entropy approximation for each of the above models.

Approximate marginals for q_x

We define the family of approximate marginal densities

$$\mathcal{Q}_x = \{ \{q_{t,t+1}^x\}_t, \{q_t^{x,f}\}_t, \{q_{t,j}^{x,g}\}_{t,j}, \{q_{t,j}^g\}_{t,j} \}. \quad (10)$$

where we assign the density $q_{t,t+1}^x(\mathbf{x}_t, \mathbf{x}_{t+1})$ to $\Psi_{t,t+1}(\mathbf{x}_t, \mathbf{x}_{t+1})$ and $q_{t+1,j}^x(\mathbf{x}_{t+1})$ to $\phi_{t+1,j}(\mathbf{x}_{t+1})$. This family can be viewed as a set of marginals that defines a joint density

$$q^x(\mathbf{X}) \propto \frac{\prod_t q_{t,t+1}^x(\mathbf{x}_t, \mathbf{x}_{t+1})}{\prod_t q_{t+1}^{x,f}(\mathbf{x}_{t+1})} \times \prod_{t,j} \frac{q_{t+1,j}^x(\mathbf{x}_{t+1})}{q_{t+1,j}^{x,l}(\mathbf{x}_{t+1})}.$$

As we will see later, this representation of q_x is exact up to the normalisation constant, however, the factors themselves are only approximations of the marginal densities. and we assume that expectation constraints w.r.t. $\mathbf{f}(\mathbf{x}_{t+1})$ and $\mathbf{l}_j(\mathbf{x}_{t+1})$ hold between $q_{t,t+1}^x$ and $q_{t+1,t+2}^x$ and $q_{t,t+1}^x$ and $q_{t+1,j}^x$ respectively. The families \mathbf{f} and \mathbf{l}_j will be restricted Gaussians, that is, \mathbf{f} and \mathbf{l}_j will correspond to a sparse Gaussian Markov random fields. In many cases we consider $\phi_{t+1,j}$ s that depend on a subset I_j of \mathbf{x}_{t+1} s elements or only one element x_{t+1}^j , that is, we have $\phi_{t+1,j}(\mathbf{x}_{t+1}; \mathbf{y}_{t+1}) = \phi_{t+1,j}(x_{t+1}^{I_j}; \mathbf{y}_{t+1})$ or $\phi_{t+1,j}(\mathbf{x}_{t+1}; \mathbf{y}_{t+1}) = \phi_{t+1,j}(x_{t+1}^j; \mathbf{y}_{t+1}^j)$ and thus $\mathbf{l}_j(\mathbf{x}_{t+1}) = \mathbf{l}_j(x_{t+1}^{I_j})$ or $\mathbf{l}_j(\mathbf{x}_{t+1}) = \mathbf{l}_j(x_{t+1}^j)$. The choice of these sufficient statistics depends on the specific model at hand and, as we will see later, our selection criteria will be the tractability of the computations they require, that is, the sparsity they generate.

Since we ...

The expectation constraints can be written as

$$\langle \mathbf{f}(\mathbf{x}_{t+1}) \rangle_{q_{t,t+1}^x} = \langle \mathbf{f}(\mathbf{x}_{t+1}) \rangle_{q_{t+1}^{x,f}} \quad \text{and} \quad \langle \mathbf{f}(\mathbf{x}_{t+1}) \rangle_{q_{t+1,t+2}^x} = \langle \mathbf{f}(\mathbf{x}_{t+1}) \rangle_{q_{t+1}^{x,f}}, \quad (11)$$

and

$$\langle \mathbf{l}_j(\mathbf{x}_{t+1}) \rangle_{q_{t,t+1}^x} = \langle \mathbf{l}_j(\mathbf{x}_{t+1}) \rangle_{q_{t+1,j}^{x,l}} \quad \text{and} \quad \langle \mathbf{l}_j(\mathbf{x}_{t+1}) \rangle_{q_{t+1,j}^x} = \langle \mathbf{l}_j(\mathbf{x}_{t+1}) \rangle_{q_{t+1,j}^{x,l}}, \quad (12)$$

respectively.

As a consequence we introduce the approximate entropy

$$-\tilde{H}(\mathcal{Q}_x) = \sum_t \langle \log q_{t,t+1}^x \rangle_{q_{t,t+1}^x} - \sum_t \langle \log q_t^{x,f} \rangle_{q_t^{x,f}} + \sum_{t,j} \langle \log q_{t,j}^x \rangle_{q_{t,j}^x} - \sum_{t,j} \langle \log q_{t,j}^{x,l} \rangle_{q_{t,j}^{x,l}} \quad (13)$$

and Lagrange multiplier terms

$$\begin{aligned} C(\mathcal{Q}_x, \Lambda_x) = & \sum_t \lambda_{t+1}^\beta \cdot \left[\langle \mathbf{f}(\mathbf{x}_{t+1}) \rangle_{q_{t+1}^{x,f}} - \langle \mathbf{f}(\mathbf{x}_{t+1}) \rangle_{q_{t,t+1}^x} \right] + \sum_t \lambda_{t+1}^\alpha \cdot \left[\langle \mathbf{f}(\mathbf{x}_{t+1}) \rangle_{q_{t+1}^{x,f}} - \langle \mathbf{f}(\mathbf{x}_{t+1}) \rangle_{q_{t+1,t+2}^x} \right] \\ & \sum_{t,j} \lambda_{t+1,j}^0 \cdot \left[\langle \mathbf{l}(\mathbf{x}_{t+1}) \rangle_{q_{t+1,j}^{x,l}} - \langle \mathbf{l}(\mathbf{x}_{t+1}) \rangle_{q_{t,t+1}^x} \right] + \sum_{t,j} \lambda_{t+1,j}^l \cdot \left[\langle \mathbf{l}(\mathbf{x}_{t+1}) \rangle_{q_{t+1,j}^{x,l}} - \langle \mathbf{l}(\mathbf{x}_{t+1}) \rangle_{q_{t+1,j}^x} \right] \end{aligned} \quad (14)$$

Approximate marginals for q_A

In case of q_A we define the family

$$\mathcal{Q}_A = \{q_0^A, \{q_{ij}^A\}_{ij}, \{q_{ij}^{A,g}\}_{ij}\} \quad (15)$$

where we associate q_0^A with the exponentiated quadratic form in q_A and we associate $q_{ij}^A(a_{ij})$ with the prior $p_0(a_{ij})$. We also define a set of expectation constraints between the members of \mathcal{Q}_A . Because of latent Gaussian nature of q_A , we will require expectation constraints up to second order. That is, we choose $\mathbf{g}(z) = (z, -z^2/2)$ and we set the constraints $\langle \mathbf{g}(a_{ij}) \rangle_{q_{ij}^A} = \langle \mathbf{g}(a_{ij}) \rangle_{q_0^A}$ and $\langle \mathbf{g}(a_{ij}) \rangle_{q_{ij}^{A,g}} = \langle \mathbf{g}(a_{ij}) \rangle_{q_0^A}$. The family \mathcal{Q}_A defines a set of marginal densities that corresponds to a density written in the form

$$q^A(\mathbf{A}) \propto q_0^A(\mathbf{A}) \prod_{ij} \frac{q_{ij}^A(a_{ij})}{q_{ij}^{A,g}(a_{ij})}. \quad (16)$$

The corresponding entropy approximation will be defined as

$$-\tilde{H}(\mathcal{Q}_A) = \langle \log q_0^A \rangle_{q_0^A} + \sum_{ij} \langle \log q_{ij}^A \rangle_{q_{ij}^A} - \sum_{ij} \langle \log q_{ij}^{A,g} \rangle_{q_{ij}^{A,g}}, \quad (17)$$

and a set of Lagrangian term for the expectation constraints can be written as

$$C(\mathcal{Q}_A, \Lambda_A) = \sum_{ij} \lambda_{g,ij}^A \cdot \left[\langle \mathbf{g}(a_{ij}) \rangle_{q_{ij}^{A,g}} - \langle \mathbf{g}(a_{ij}) \rangle_{q_{ij}^A} \right] + \sum_{ij} \lambda_{0,ij}^A \cdot \left[\langle \mathbf{g}(a_{ij}) \rangle_{q_{ij}^{A,g}} - \langle \mathbf{g}(a_{ij}) \rangle_{q_0^A} \right]. \quad (18)$$

Issues related to the sparsity constraint on \mathbf{A} will be discussed in Section (?).

Approximate marginals for q_B

In case of q_B we define the family

$$\mathcal{Q}_B = \{q_0^B, \{q_{ij}^B\}_{ij}, \{q_{ij}^{B,g}\}_{ij}\} \quad (19)$$

and we assign q_0^B with the exponentiated quadratic form in q_B and we assign $q_{ij}^B(b_{ij})$ with the prior $p_0(b_{ij})$. We define the expectation constraints similarly as in the section above. The corresponding approximating density can be viewed as

$$q^B(\mathbf{B}) \propto q_0^B(\mathbf{B}) \prod_{ij} \frac{q_{ij}^B(b_{ij})}{q_{ij}^{B,g}(b_{ij})}. \quad (20)$$

The approximate entropy term will be

$$-\tilde{H}(\mathcal{Q}_B) = \langle \log q_0^B \rangle_{q_0^B} + \sum_{ij} \langle \log q_{ij}^B \rangle_{q_{ij}^B} - \sum_{ij} \langle \log q_{ij}^{B,g} \rangle_{q_{ij}^{B,g}} \quad (21)$$

and the Langrangian term corresponding to the expectation constraints will be written as

$$C(\mathcal{Q}_B, \Lambda_B) = \sum_{ij} \lambda_{g,ij}^B \cdot \left[\langle \mathbf{g}(b_{ij}) \rangle_{q_{ij}^{B,g}} - \langle \mathbf{g}(b_{ij}) \rangle_{q_{ij}^B} \right] + \sum_{ij} \lambda_{0,ij}^B \cdot \left[\langle \mathbf{g}(b_{ij}) \rangle_{q_{ij}^{B,g}} - \langle \mathbf{g}(b_{ij}) \rangle_{q_0^B} \right]. \quad (22)$$

Approximate marginals for q_Q

We do not define any approximation for q_Q , we assume that the free form variational update can be done analytically and all the expectation and the entropy term needed are also tractable.

Free energy optimisation and message passing

The Lagrangian corresponding to the free energy minimisation expressed in term of $\mathcal{Q}_x, \mathcal{Q}_A, \mathcal{Q}_B$ and q_Q is

$$\begin{aligned} L(\mathcal{Q}_x, \mathcal{Q}_A, \mathcal{Q}_Q, \mathcal{Q}_B, \Lambda_x, \Lambda_A, \Lambda_Q, \Lambda_B) &= - \langle \log p(\mathbf{Y}, \mathbf{X}, \mathbf{A}, \mathbf{Q}, \mathbf{B} | \mathbf{U}) \rangle_{\mathcal{Q}_x, \mathcal{Q}_A, \mathcal{Q}_Q, \mathcal{Q}_B} \\ &\quad - \tilde{H}(\mathcal{Q}_x) - \tilde{H}(\mathcal{Q}_A) - H(q_Q) - \tilde{H}(\mathcal{Q}_B) \\ &\quad + C(\mathcal{Q}_x, \Lambda_x) + C(\mathcal{Q}_A, \Lambda_A) + C(\mathcal{Q}_B, \Lambda_B) + \text{normalisation constraints.} \end{aligned}$$

In addition, let us define $\text{Collapse}(p(\mathbf{z}); \mathbf{f})$ as the (moment) projection of $p(\mathbf{x})$ into the exponential family defined by $\mathbf{f}(\mathbf{x})$, that is,

$$\text{Collapse}[p(\mathbf{z}); \mathbf{f}] = \underset{\theta}{\text{argmin}} D[p(\mathbf{z}) || \exp(\theta \cdot \mathbf{f}(\mathbf{z}) - \log Z(\theta))].$$

To differentiate between the free form densities, for example q_A , and the corresponding family of marginals, for example \mathcal{Q}_A , we use a different kind of notation for expectations. The quantities corresponding to \mathbf{h}_A and \mathbf{Q}_A will be denoted by $\langle \mathbf{h}_A \rangle_{\mathcal{Q}_A}$ and $\langle \mathbf{Q}_A \rangle_{\mathcal{Q}_A}$ and the same rule applies for q_B and \mathcal{Q}_B . For q_x and \mathcal{Q}_x we use the shortcut notation $\langle \log \Psi_{t,t+1}(\mathbf{x}_t, \mathbf{x}_{t+1}) \rangle_{\mathcal{Q}_x}$, but the same usage applies.

From the stationary conditions of L w.r.t the members of the families $\mathcal{Q}_x, \mathcal{Q}_A$ and \mathcal{Q}_B we can derive the form of the approximating densities. The stationarity conditions corresponding to the expectation constraints will be used to define a message passing algorithm.

These are as follows

(1) for the members of the family \mathcal{Q}_x we have

$$q_{t,t+1}^x(\mathbf{x}_t, \mathbf{x}_{t+1}) \propto \exp \left\{ \langle \log \Psi_{t,t+1}(\mathbf{x}_t, \mathbf{x}_{t+1}) \rangle_{\mathcal{Q}_x} \right. \quad (23)$$

$$\left. + \lambda_t^\alpha \cdot \mathbf{f}(\mathbf{x}_t) + \lambda_{t+1}^\beta \cdot \mathbf{f}(\mathbf{x}_{t+1}) + \sum_j \lambda_{t+1,j}^0 \cdot \mathbf{l}_j(\mathbf{x}_{t+1}) \right\}, \quad (24)$$

$$q_{t+1,j}^{x,l}(\mathbf{x}_{t+1}) \propto \psi_{t+1,j}(\mathbf{x}_{t+1}; \mathbf{y}_{t+1}) \times \exp\{\lambda_{t+1,j}^l \cdot \mathbf{l}_j(\mathbf{x}_{t+1})\}, \quad (25)$$

$$q_{t+1,j}^{x,l}(\mathbf{x}_{t+1}) \propto \exp\{(\lambda_{t+1,j}^0 + \lambda_{t+1,j}^l) \cdot \mathbf{l}_j(\mathbf{x}_{t+1})\}, \quad (26)$$

$$q_{t+1}^{x,f}(\mathbf{x}_{t+1}) \propto \exp\{(\lambda_{t+1}^\alpha + \lambda_{t+1}^f) \cdot \mathbf{f}(\mathbf{x}_{t+1})\}, \quad (27)$$

and the corresponding expectation constraints lead the update equations

$$[\boldsymbol{\lambda}_{t+1,j}^l]^{new} = \text{Collapse}(q_{t,t+1}^x[\mathbf{x}_{t+1}]; \mathbf{l}_j) - \boldsymbol{\lambda}_{t+1,j}^0, \quad (28)$$

$$[\boldsymbol{\lambda}_{t+1,j}^0]^{new} = \text{Collapse}(q_{t+1,j}^x[\mathbf{x}_{t+1}]; \mathbf{l}_j) - \boldsymbol{\lambda}_{t+1,j}^l, \quad (29)$$

$$[\boldsymbol{\lambda}_{t+1}^\alpha]^{new} = \text{Collapse}(q_{t,t+1}^x[\mathbf{x}_{t+1}]; \mathbf{f}) - \boldsymbol{\lambda}_{t+1}^\beta, \quad (30)$$

$$[\boldsymbol{\beta}_{t+1}^\alpha]^{new} = \text{Collapse}(q_{t+1,t+2}^x[\mathbf{x}_{t+1}]; \mathbf{f}) - \boldsymbol{\lambda}_{t+1}^\alpha, \quad (31)$$

(2) for the members of the family \mathcal{Q}_A we have

$$q_0^A \propto \exp \left\{ [\mathbf{A}]_c^T \langle \mathbf{h}_A \rangle_{\mathcal{Q}_A} - \frac{1}{2} [\mathbf{A}]_c^T \langle \mathbf{Q}_A \rangle_{\mathcal{Q}_A} [\mathbf{A}]_c + \sum_{ij} \boldsymbol{\lambda}_{0,ij}^A \cdot \mathbf{g}(a_{ij}) \right\}, \quad (32)$$

$$q_{ij}^A \propto p_0(a_{ij}) \times \exp \left\{ \boldsymbol{\lambda}_{g,ij}^A \cdot \mathbf{g}(a_{ij}) \right\}, \quad (33)$$

$$q_{ij}^{A,g} \propto \exp \left\{ (\boldsymbol{\lambda}_{0,ij}^A + \boldsymbol{\lambda}_{g,ij}^A) \cdot \mathbf{g}(a_{ij}) \right\}, \quad (34)$$

and the corresponding expectation constraints lead the equations

$$[\boldsymbol{\lambda}_{0,ij}^A]^{new} = \text{Collapse}[q_{ij}^A(a_{ij}); \mathbf{g}] - \boldsymbol{\lambda}_{g,ij}^A, \quad (35)$$

$$[\boldsymbol{\lambda}_{g,ij}^A]^{new} = \text{Collapse}[q_0^A(a_{ij}); \mathbf{g}] - \boldsymbol{\lambda}_{0,ij}^A, \quad (36)$$

(3) for the members of the family \mathcal{Q}_B we have

$$q_0^B \propto \exp \left\{ [\mathbf{B}]_c^T \langle \mathbf{h}_B \rangle_{\mathcal{Q}_B} - \frac{1}{2} [\mathbf{B}]_c^T \langle \mathbf{Q}_B \rangle_{\mathcal{Q}_B} [\mathbf{B}]_c + \sum_{ij} \boldsymbol{\lambda}_{0,ij}^B \cdot \mathbf{g}(b_{ij}) \right\}, \quad (37)$$

$$q_{ij}^B \propto p_0(b_{ij}) \times \exp \left\{ \boldsymbol{\lambda}_{g,ij}^B \cdot \mathbf{g}(b_{ij}) \right\}, \quad (38)$$

$$q_{ij}^{B,g} \propto \exp \left\{ (\boldsymbol{\lambda}_{0,ij}^B + \boldsymbol{\lambda}_{g,ij}^B) \cdot \mathbf{g}(b_{ij}) \right\}, \quad (39)$$

and the corresponding expectation constraints lead the equations

$$[\boldsymbol{\lambda}_{0,ij}^B]^{new} = \text{Collapse}[q_{ij}^B(a_{ij}); \mathbf{g}] - \boldsymbol{\lambda}_{g,ij}^B, \quad (40)$$

$$[\boldsymbol{\lambda}_{g,ij}^B]^{new} = \text{Collapse}[q_0^B(a_{ij}); \mathbf{g}] - \boldsymbol{\lambda}_{0,ij}^B. \quad (41)$$

Details of the collapse operations

The collapse operations on $q_{t,t+1}^x$

Let $\boldsymbol{\lambda}_t^\alpha = (\mathbf{h}_t^\alpha, \mathbf{Q}_t^\alpha)$, $\boldsymbol{\lambda}_t^\beta = (\mathbf{h}_t^\beta, \mathbf{Q}_t^\beta)$ and $\boldsymbol{\lambda}_{t+1,j}^0 = (\mathbf{h}_{t+1,j}^0, \mathbf{Q}_{t+1,j}^0)$. Then the density $q_{t,t+1}^x$ is a multivariate Gaussian with canonical parameters

$$\tilde{\mathbf{h}}_{t,t+1}^x = \langle \mathbf{h}_{t,t+1}^x \rangle_{\mathcal{Q}_x} + \begin{bmatrix} \mathbf{h}_{t+1,j}^\alpha \\ \mathbf{h}_{t+1,j}^\beta + \sum_j \mathbf{h}_{t+1,j}^0 \end{bmatrix}, \quad (42)$$

$$= \begin{bmatrix} -\langle \mathbf{A} \rangle_{\mathcal{Q}_A}^T \langle \mathbf{Q} \rangle_{q_Q} \langle \mathbf{B} \rangle_{\mathcal{Q}_B} \mathbf{u}_t + \mathbf{h}_{t+1,j}^\alpha \\ \mathbf{h}_{t+1,j}^\beta + \sum_j \mathbf{h}_{t+1,j}^0 \end{bmatrix}, \quad (43)$$

$$\tilde{\mathbf{Q}}_{t,t+1}^x = \langle \mathbf{Q}_{t,t+1}^x \rangle_{\mathcal{Q}_x} + \begin{bmatrix} \mathbf{Q}_t^\alpha & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_{t+1}^\beta + \sum_j \mathbf{Q}_{t+1,j}^0 \end{bmatrix} \quad (44)$$

$$= \begin{bmatrix} \langle \mathbf{A}^T \langle \mathbf{Q} \rangle_{q_Q} \mathbf{A} \rangle_{\mathcal{Q}_A} + \mathbf{Q}_t^\alpha & -\langle \mathbf{A} \rangle_{\mathcal{Q}_A}^T \langle \mathbf{Q} \rangle_{q_Q} \\ -\langle \mathbf{Q} \rangle_{q_Q} \langle \mathbf{A} \rangle_{\mathcal{Q}_A} & \langle \mathbf{Q} \rangle_{q_Q} + \mathbf{Q}_{t+1}^\beta + \sum_j \mathbf{Q}_{t+1,j}^0 \end{bmatrix} \quad (45)$$

and the collapse operations are equivalent to computing the mean and covariance values corresponding to the non-zeros in \mathbf{Q}_t^α , \mathbf{Q}_t^β and $\mathbf{Q}_{t+1,j}^0$. Moreover, since further covariance values are needed for \mathcal{Q}_A and q_Q , it is necessary to compute all covariance values corresponding to the non-zeros in $\langle \mathbf{Q}_{t,t+1}^x \rangle_{\mathcal{Q}_x}$. This can be achieved by applying sparse (reordered) Cholesky factorisation and partial matrix inversion that makes use of the Cholesky factor to solve the Takahashi equations. As expected all collapse operations on $q_{t,t+1}^x$ are connected, for this reason, we will detail each of them and then we assess how can we perform them in the most efficient way.

The Collapse $[q_{t,t+1}^x(\mathbf{x}_{t+1}); \mathbf{f}]$ operation. This operation involves projecting the the marginal density $q_{t,t+1}^x$ into the (sparse) Gaussian family defined by \mathbf{f} . The optimisation corresponding to this collapse step can be written as

$$(46)$$

Sparse approximations with known structure for \mathbf{A} and \mathbf{Q}

In this section we discuss the sparsity issue that are essential for the scalability of the procedure we present. The basic assumptions we started with are that both \mathbf{A} and \mathbf{Q} are sparse. We also assume that \mathbf{B} is a full matrix, although if we want to assume some structural sparsity, the issue can be addressed in the same way as for \mathbf{A} .

An crucial observation is that in all collapse steps we perform, we act on a sparse precision matrix and we compute all correlation values corresponding to the non-zero positions in this precision matrix. We argue that these correlation values are sufficient to keep the algorithm running, that is, there is no need to compute extra correlations and the correlations we compute always lead to positive semi-definite precision matrices. Given our observation about the computation of the correlation values, the first argument is quite intuitive: wherever there is a term with quadratic interaction between two variables of the same type be it \mathbf{X} , \mathbf{A} , \mathbf{B} the moments of the other variables be it quadratic or just means are computed by their corresponding inference scheme and thus this property is conserved. The second argument related to positive (semi)-definiteness requires a closer look.

First we explore the connection between \mathbf{X} and \mathbf{A} . We have $\langle \mathbf{Q}_{t,t+1}^x \rangle_{\mathcal{Q}_x} = \langle \mathbf{Q}_{t,t+1}^x \rangle_{q_0^A q_Q}$. Since $(\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t)^T \mathbf{Q} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t) \geq 0$ for any $\mathbf{x}_t, \mathbf{x}_{t+1}, \mathbf{A}$ and positive semi-definite \mathbf{Q} , it follows that $\langle (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t)^T \langle \mathbf{Q} \rangle_{q_Q} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t) \rangle_{q_0^A} \geq 0$ for any positive semi-definite $\langle \mathbf{Q} \rangle_{q_Q}$ thus $\langle \mathbf{Q}_{t,t+1}^x \rangle_{\mathcal{Q}_x}$ is positive semi-definite for all semi-definite $\langle \mathbf{Q} \rangle_{q_Q}$. The argument for $\langle \mathbf{Q}_A \rangle_{\mathcal{Q}_A}$ is as follows: first we construct the matrix $\sum_t \langle \mathbf{x}_t \mathbf{x}_t^T \rangle_{\mathcal{Q}_x} \otimes \langle \mathbf{Q} \rangle_{q_Q} = \sum_t \langle \mathbf{x}_t \mathbf{x}_t^T \rangle_{q_{t,t+1}^x} \otimes \langle \mathbf{Q} \rangle_{q_Q}$ which, since it is a Kronecker product of positive semi-definite matrices, it is also positive semi-definite. Now, by eliminating the the rows and columns corresponding to the a-priori set structural zeros in $[\mathbf{A}]_c$ we arrive to a principal minor that also has to be positive semi-definite. The (i, j) block of $\langle \mathbf{Q}_A \rangle_{\mathcal{Q}_A}$ is $\langle x_t^i x_t^j \rangle [\langle \mathbf{Q} \rangle_{q_Q}]_{I_i, I_j}$ where I_i and I_j are the supports of \mathbf{a}_i and \mathbf{a}_j respectively. Since the elements of $\mathbf{A}^T \mathbf{Q} \mathbf{A}$ are $\mathbf{a}_{I_i, i}^T \mathbf{Q}_{I_i, I_j} \mathbf{a}_{I_j, j}$, it follows that $\langle x_t^i x_t^j \rangle$ is always computed when needed, that is, when $I_i \cap I_j \neq \emptyset$. This shows that the sparsity structures can be exploited and that the resulting inference keeps the matrices positive semi-definite.

A Appendix

A.1 Structured mean field approximation

$$D[q_1(\mathbf{x}_{I_1})q_2(\mathbf{x}_{I_2})\dots q(\mathbf{x}_{I_K})||p(\mathbf{x})] = -\langle \log p(\mathbf{x}) \rangle_{q_1q_2\dots q_K} + \sum_j \langle \log q_j(x_{I_j}) \rangle_{q_j} \quad (47)$$

$$q_j(\mathbf{x}_{I_j}) \propto \exp\{\langle \log p(\mathbf{x}) \rangle_{\setminus q_j}\} \quad (48)$$

A.2 Inference using expectation constraints

Let us assume that we have a density that factors as

$$p(\mathbf{x}) = \frac{1}{Z_p} \prod_j \Psi_j(\mathbf{x}_{I_j})$$

and let us assume that the factors Ψ_j represent this density at the largest resolution, that is, no Ψ_j can be factored further. Now, say, we want to approximate p with a density q that has some desirable features, say, easy access to marginals or expectations (features common in exponential families), that p has not. The Kullback-Leibler divergence $D[\cdot||\cdot]$ or some of its versions like the α -divergence come as natural measures to minimise. The $D[p||q]$ version is typically out of question because it typically requires computations of expectations and marginals that we want to get around of in the first place. A few examples are: (1) say, we want to do a mean field approximation in the $D[p(\mathbf{x})||\prod_j q(x_i)]$ sense, then, the optimality conditions lead to $q(x_i) = p(x_i)$ or (2) we want to q to have an exponential form $q(\mathbf{x}; \boldsymbol{\theta}) = \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}) - \log Z(\boldsymbol{\theta}))$, then we have to compute $\langle \mathbf{f}_x \rangle_p$ followed by the moment transformation specific to this family. All these alternatives, involve the very steps we want to avoid.

However, the other version $D[q||p]$ requires the computation of the entropy of q which can become complicated. Nonetheless, as a trade off we get a relatively easy access to p and we can exploit its factorisation. The latter, so called variational version of $D[\cdot||\cdot]$ reads as

$$D[q||p] = \langle \log q(\mathbf{x}) \rangle_{q(\mathbf{x})} - \sum_j \langle \log \Psi_j(\mathbf{x}_{I_j}) \rangle_{q(\mathbf{x}_{I_j})} + \log Z_p \quad (49)$$

In the above mentioned example we obtain the objectives

$$D[q||p] = \sum_j \langle \log q_j(x_j) \rangle_{q_j(x_j)} - \sum_j \langle \log \Psi_j(\mathbf{x}_{I_j}) \rangle_{\prod_{k \in I_j} q_k(x_k)} \quad (50)$$

and

$$D[q||p] = \boldsymbol{\theta} \cdot \partial_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta}) - \log Z(\boldsymbol{\theta}) - \sum_j \langle \log \Psi_j(\mathbf{x}_{I_j}) \rangle_{q(\mathbf{x}_{I_j})} \quad (51)$$

which are more manageable than the former version $D[p||q]$.

In the following we'll try to go beyond the mean field approach in Equation (50). Let us assume that we have constructed a junction tree on the graph defined by the interactions in the cliques I_j , and let us denote the node sets of this tree by C_1, \dots, C_K and let us use the notation S_1, \dots, S_L for the separator corresponding to this tree. Moreover, let us restructure the factors $\Psi_j(\mathbf{x}_{I_j})$ by forming new factors $\Phi_j(\mathbf{x}_{C_j})$ corresponding to the clique structure of the junction tree. A factor $\Psi_j(\mathbf{x}_{I_j})$ can be attached to any factor $\Phi_j(\mathbf{x}_{C_j})$ where $I_k \subseteq C_j$ holds. We will assume that the approximating density follows the structure of the junction tree and thus can be written as

$$q(\mathbf{x}) = \frac{\prod_j q(\mathbf{x}_{C_j})}{\prod_k q(\mathbf{x}_{S_k})^{n_k-1}}.$$

where n_k is the number of cliques a separator appears in. Note that since S_k is separator $n_k \geq 2$. Since q has a clique-tree structure we can write the corresponding variational KL-divergence as

$$D[q||p] = \sum_j \langle \log q(\mathbf{x}_{C_j}) \rangle_{q(\mathbf{x}_{C_j})} - \sum_k (n_k - 1) \langle \log q(\mathbf{x}_{S_k}) \rangle_{q(\mathbf{x}_{S_k})} - \sum_j \langle \log \Phi_j(\mathbf{x}_{C_j}) \rangle_{q(\mathbf{x}_{C_j})} + \log Z_p.$$

Clearly, in order to turn this into an optimisation problem we have to replace q by a family of approximate marginals $\mathcal{Q} = \{\{q_{C_j}(\mathbf{x}_{C_j})\}_j, \{q_{S_k}(\mathbf{x}_{S_k})\}_k\}$ where we assume that the constraints $\int d\mathbf{x}_{C_j \setminus S_k} q_{C_j}(\mathbf{x}_{C_j}) = q_{S_k}(\mathbf{x}_{S_k})$ hold for any $S_k \subseteq C_j$. With these constraints we define the variational objective

$$F(\mathcal{Q}) \equiv \sum_j \langle \log q_{C_j}(\mathbf{x}_{C_j}) \rangle_{q_{C_j}(\mathbf{x}_{C_j})} - \sum_k (n_k - 1) \langle \log q_{S_k}(\mathbf{x}_{S_k}) \rangle_{q_{S_k}(\mathbf{x}_{S_k})} - \sum_j \langle \log \Phi_j(\mathbf{x}_{C_j}) \rangle_{q_{C_j}(\mathbf{x}_{C_j})} \quad (52)$$

and the Lagrangian

$$L(\mathcal{Q}, \Lambda) \equiv F(\mathcal{Q}) + \sum_{j,k:S_k \subseteq C_j} \int d\mathbf{x}_{S_k} \mu_{j,k}(\mathbf{x}_{S_k}) [q_{S_k}(\mathbf{x}_{S_k}) - q_{C_j}(\mathbf{x}_{S_k})] \\ + \sum_j z_{C_j} [\int d\mathbf{x}_{C_j} q_{C_j}(\mathbf{x}_{C_j}) - 1] + \sum_k z_{S_k} [1 - \int d\mathbf{x}_{S_k} q_{S_k}(\mathbf{x}_{S_k})].$$

From the stationarity conditions we obtain that

$$q_{C_j}(\mathbf{x}_{C_j}) = \Phi_j(\mathbf{x}_{C_j}) \times \exp \left\{ \sum_{k:S_k \subseteq C_j} \mu_{j,k}(\mathbf{x}_{S_k}) - z_{C_j} \right\} \\ q_{S_k}(\mathbf{x}_{S_k}) = \exp \left\{ \frac{1}{n_k - 1} \sum_{j:C_j \supseteq S_k} \mu_{j,k}(\mathbf{x}_{S_k}) - z_{S_k} / (n_k - 1) \right\} \\ q_{S_k}(\mathbf{x}_{S_k}) = q_{C_j}(\mathbf{x}_{S_k}) \quad \text{for any } C_j \cap S_k \neq \emptyset.$$

The latter equation can be written as

$$\exp \left\{ \frac{1}{n_k - 1} \sum_{j':C_j \supseteq S_k} \mu_{j',k}(\mathbf{x}_{S_k}) - z_{S_k} / (n_k - 1) \right\} = \int d\mathbf{x}_{C_j \setminus S_k} \Phi_j(\mathbf{x}_{C_j}) \times \exp \left\{ \sum_{k':S_k \subseteq C_j} \mu_{j,k'}(\mathbf{x}_{S_{k'}}) - z_{C_j} \right\} \quad (53)$$

We introduce a non-singular linear transformation of the Lagrange multipliers

$$\mu_{jk}(\mathbf{x}_{S_k}) = \sum_{j':C_{j'} \supseteq S_k} \hat{\mu}_{j'k}(\mathbf{x}_{S_k}).$$

By rearranging the terms we get

$$\exp \left\{ \mu_{jk}(\mathbf{x}_{S_k}) + \hat{\mu}_{jk}(\mathbf{x}_{S_k}) \right\} \propto \int d\mathbf{x}_{C_j \setminus S_k} \Phi_j(\mathbf{x}_{C_j}) \times \exp \left\{ \sum_{k':S_{k'} \subseteq C_j} \mu_{jk'}(\mathbf{x}_{S_{k'}}) \right\} \\ \propto \int d\mathbf{x}_{C_j \setminus S_k} \Phi_j(\mathbf{x}_{C_j}) \times \exp \left\{ \sum_{k':S_{k'} \subseteq C_j} \sum_{j' \neq j: C_{j'} \supseteq S_k} \hat{\mu}_{j'k'}(\mathbf{x}_{S_{k'}}) \right\}.$$

Using these notations and stationarity conditions one can derive the fixed point equations/updates known as message passing. This has the form

$$q_{C_j}(\mathbf{x}_{C_j}) \propto \Phi_j(\mathbf{x}_{C_j}) \times \exp \left\{ \sum_{k':S_{k'} \subseteq C_j} \mu_{jk'}(\mathbf{x}_{S_{k'}}) \right\}, \quad (54)$$

$$q_{S_k}(\mathbf{x}_{S_k}) \propto \exp \left\{ \sum_{j':C_{j'} \supseteq S_k} \hat{\mu}_{j'k}(\mathbf{x}_{S_k}) \right\}, \quad (55)$$

$$\hat{\mu}_{jk}(\mathbf{x}_{S_k}) \stackrel{\circ}{=} \log \left\{ \int d\mathbf{x}_{C_j \setminus S_k} q_{C_j}(\mathbf{x}_{C_j}) \right\} - \mu_{jk}(\mathbf{x}_{S_k}), \quad (56)$$

$$\mu_{jk}(\mathbf{x}_{S_k}) \stackrel{\circ}{=} \log \{ q_{S_k}(\mathbf{x}_{S_k}) \} - \hat{\mu}_{jk}(\mathbf{x}_{S_k}). \quad (57)$$

where we use the symbol “ \doteq ” to denote equality up to a constant, for example, $2x \doteq 2x + 1$. The functions μ_{jk} and $\hat{\mu}_{jk}$ are called messages. There are several specialisations of this algorithm, in the following we detail two of these: (1) when the constraints $q_{S_k}(\mathbf{x}_{S_k}) = q_{C_j}(\mathbf{x}_{S_k})$ is replaced by moment constraints and (2) when we no longer operate on clique-trees that is C_j and S_j are arbitrary.

A.2.1 Expectation constraints

In many cases the integrals in equations (53) and (61) are not analytically computable and having free form messages as shown above can be computationally demanding. To get a handle on this issue a plausible approach is to parameterise the messages (Lagrange multipliers). However, in many cases the marginal consistency constraints in (53) cannot be satisfied and further approximations have to be applied. First we show that parameterised Lagrange multipliers or messages are formally equivalent to moment matching constraints (instead of marginal matching constraints). Say, we want all messages to be in the linear parametric family defined by $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{z})$. Then the Lagrange multiplier terms have the form

$$\int d\mathbf{x}_{S_k} [\boldsymbol{\mu}_{jk} \cdot \mathbf{f}(\mathbf{x}_{S_k})] [q_{S_k}(\mathbf{x}_{S_k}) - q_{C_j}(\mathbf{x}_{S_k})] = \boldsymbol{\mu}_{jk} \cdot \left[\int d\mathbf{x}_{S_k} q_{S_k}(\mathbf{x}_{S_k}) \mathbf{f}(\mathbf{x}_{S_k}) - \int d\mathbf{x}_{S_k} q_{C_j}(\mathbf{x}_{S_k}) \mathbf{f}(\mathbf{x}_{S_k}) \right] \quad (58)$$

an thus this can be formally interpreted as replacing the marginal matching constraints $q_{C_j}(\mathbf{x}_{S_k}) = q_{S_k}(\mathbf{x}_{S_k})$ with the moment matching constraints $\langle \mathbf{f}(\mathbf{x}_{S_k}) \rangle_{q_{C_j}} = \langle \mathbf{f}(\mathbf{x}_{S_k}) \rangle_{q_{S_k}}$. This implies that the the marginal matching in Equation (53) is replaced my moment matching. Let us define

$$\text{Collapse}[p(\mathbf{z}); \mathbf{f}] = \underset{\boldsymbol{\theta}}{\text{argmin}} D[p(\mathbf{z}) || \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{z}) - \log Z(\boldsymbol{\theta}))]$$

and assume that it is unique. Then the message passing algorithm can be rewritten as

$$q_{C_j}(\mathbf{x}_{C_j}) \propto \Phi_j(\mathbf{x}_{C_j}) \times \exp \left\{ \sum_{k': S_{k'} \subseteq C_j} \boldsymbol{\mu}_{jk'} \cdot \mathbf{f}(\mathbf{x}_{S_{k'}}) \right\}, \quad (59)$$

$$q_{S_k}(\mathbf{x}_{S_k}) \propto \exp \left\{ \sum_{j': C_{j'} \supset S_j} \hat{\boldsymbol{\mu}}_{j'k} \cdot \mathbf{f}(\mathbf{x}_{S_k}) \right\}, \quad (60)$$

$$\hat{\boldsymbol{\mu}}_{jk} = \text{Collapse}[q_{C_j}(\mathbf{x}_{S_k}); \mathbf{f}] - \boldsymbol{\mu}_{jk}, \quad (61)$$

$$\boldsymbol{\mu}_{jk} = \text{Collapse}[q_{S_k}(\mathbf{x}_{S_k}); \mathbf{f}] - \hat{\boldsymbol{\mu}}_{jk}. \quad (62)$$

Note that $q_{S_k}(\mathbf{x}_{S_k})$ is already in the family defined by \mathbf{f} , therefore, the only role of $\text{Collapse}[q_{S_k}(\mathbf{x}_{S_k}); \mathbf{f}]$ is parameter identification. It is important to mention that at no point the messages $\boldsymbol{\mu}_{jk}$ and $\hat{\boldsymbol{\mu}}_{jk}$ have to represent normalizable functions, that is, valid densities. It is sufficient to have q_{S_k} and q_{C_j} normalizable. Moreover, \mathbf{f} can vary from constraint to constraint as long as manage to keep the computations under control.

A.2.2 Expectation propagation

A further algorithm follows by removing the clique-tree requirement, in this case we consider $C_j \cap C_j \not\subseteq C_j$ for any i and j and define the sets S_k as atomic sets satisfying the condition $S_k \cap C_j$ implies $S_k \subseteq C_j$. In this case we approximate the entropy by

$$\tilde{H}(\mathcal{Q}) = \sum_j \langle \log q(\mathbf{x}_{C_j}) \rangle_{q(\mathbf{x}_{C_j})} - \sum_k (n_k - 1) \langle \log q(\mathbf{x}_{S_k}) \rangle_{q(\mathbf{x}_{S_k})},$$

where n_k is again the number of sets of type C_j that S_k appears in, that is, $n_k = |\{C_j : S_k \subseteq S_j\}|$. This approximation of the entropy is similar to the Bethe entropy approximation in (?). The resulting algorithm is known as expectation propagation (?).