

Message Passing and Expectation Propagation

Efficient Inference in large scale machine learning

Christoph Dehner
Department of Informatics
Technische Universität München
dehner@in.tum.de

ABSTRACT

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

1. INTRODUCTION

Probabilistic graphical models like Markov Random Fields or Bayesian Networks provide clear and illustrative ways to describe probabilistic processes. In such models, nodes represent random variables and edges their conditional dependencies. The joint distribution of all involved variables can be expressed as a product of factors, which are observed or given specific values during modelling. As an example, a Bayesian network with three variables is given in figure 1. It defines the values of x_1 and x_3 to be conditionally independent given x_2 .

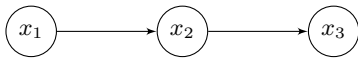


Figure 1: Bayesian network with three variables. x_1 and x_3 are conditionally independent given x_2 .

Typical tasks in such a scenario are to do Bayesian inference to extract marginal distribution of specific variables or find maximum a priori estimates. However, with an increasing number of variables these task can become computationally very expensive.

The naive approach for marginalisation is to sum out all but one variable from the joint distribution $p(X)$. This is formally shown in equation 1, where $X \setminus x_i$ describes the set of all random variables except x_i . This strategy requires exponentially in the number of variables many evaluations of the joint distribution and is thus not applicable to bigger

networks.

$$p(x_i) = \sum_{X \setminus x_i} p(\mathbf{x}) \quad (1)$$

Therefore, this paper explains more efficient algorithms to do large scale Bayesian inference in graphical models. The next chapter *todo* deals with exact inference and presents two message passing algorithms to calculate marginals and maximum a priori estimates for graphical models. Subsequently, in chapter *todo* expectation propagation as a method of approximate inference is explained.

2. EXACT INFERENCE

The intuition of the algorithms presented in this section is to exploit independence properties of the random variables. Graphical models define the joint probability to be expressed as a product of factors, considering the conditional independences expressed by the graph. For Bayesian networks those factors correspond to conditional distributions, in Markov Random Fields to clique potentials. Assuming appropriately normalized factors, the joint distribution can be written as a product according to equation 2. The index s iterates here over all factors of the graph; X_s defines the subset of variables, the factor s depends on.

$$p(X) = \prod_s f_s(X_s) \quad (2)$$

Expressing the joint distribution as the product defined by the graph allows to exchange summations and multiplications during marginalization. By this way, marginalization can often be done a lot more efficient. Equation 3 demonstrates this transformation with the Bayesian network from figure 1, whose joint distribution $p(X)$ is defined as $p(X) = p(x_1)p(x_2|x_1)p(x_3|x_1)$.

$$\begin{aligned} p(x_2) &= \sum_{x_1} \sum_{x_3} p(x_1)p(x_2|x_1)p(x_3|x_2) \\ &= \sum_{x_1} \left[p(x_1)p(x_2|x_1) \sum_{x_3} \left[p(x_3|x_2) \right] \right] \\ &= \underbrace{\left[\sum_{x_1} p(x_1)p(x_2|x_1) \right]}_{\mu_{x_1 \rightarrow x_2}} \cdot \underbrace{\left[\sum_{x_3} p(x_3|x_2) \right]}_{\mu_{x_3 \rightarrow x_2}} \end{aligned} \quad (3)$$

Here, the sum of products from the naive approach for marginalization is transformed to a product of sums, reducing the exponential computational complexity in the number of random variables to linear effort.

The brackets in the last line of equation 3 reveal a powerful interpretation of the marginalizing. The marginal distribution $p(x_2)$ consists of two messages $\mu_{x_1 \rightarrow x_2}$ and $\mu_{x_3 \rightarrow x_2}$ from its neighboring cells x_1 and x_3 . For a longer chain of random variables these messages would again be comprised by messages from their neighboring cell. Applied to a chain of arbitrary length, this method gives a recursive pattern in which messages for the marginalization are sent to the marginalized variable node from both ends of the chain.

This message passing idea is the foundation for the two inference algorithms presented in this chapter subsequently. Before that the next section introduces factor graphs, the structure these algorithms operate on.

2.1 Factor graphs

Factor graphs make the factorisation of a joint probability distribution explicit. They can be generated from Bayesian networks as well as from Markov random fields and thus allow to define inference algorithms independent of how the underlying probabilistic graphical model was introduced.

Additional to the variable nodes, a factor graph also consists of factor nodes representing the factors of the decomposed joint probability distribution as in equation 2. Edges connect the factor nodes to all variables they depend on. By this way they are a bipartite graph with variable nodes (usually visualised by circles) on the one side and factors (depicted as rectangles) on the other side.

Depending on the exact factorisation, a distribution defined by a Bayesian network or Markov random field can be represented by different factor graphs. Figure *todo* depicts a factor graph of the Bayesian network from figure 1, in which all conditional distributions are represented by separate factors. The inference algorithms on factor graphs of the following sections are valid for trees. Such factor graphs with exactly one path between any pair of nodes can be generated from undirected trees in the case of a Markov random field model as well as from directed trees and polytrees, if the factor graph is derived from a Bayesian network. More detailed description how to derive factor graphs from probabilistic graph models can be found in *todo: add reference*.

2.2 Marginalization

The sum-product algorithm presented in this section generalises the idea of exchanging summation and multiplication during marginalisation. It consists of local messages passed between factor and variable nodes in a factor graph.

As a starting point let us consider the calculation of the marginal distribution $p(x_i)$ in the factor graph of figure 2. Furthermore we define $F_s(x_i, X_s)$ for all neighbours s of x_i as the product of all factors in the respective sub-tree of the neighbour. This is visualised by blue circles in figure 2.

Analogous to the general factorisation defined by the graphical model from equation 2 the joint probability $p(\mathbf{x})$ can then be expressed as a product of those neighbour-factor of x_i :

$$p(X) = \prod_{s \in ne(x_i)} F_s(x_i, X_s) \quad (4)$$

Inserting this representation of the joint distribution in the naive marginalisation formula from equation 1 allows to exchange sum and product. By this way the marginal distribution of x_i is expressed as the product of messages $\mu_{f_s \rightarrow x_i}$

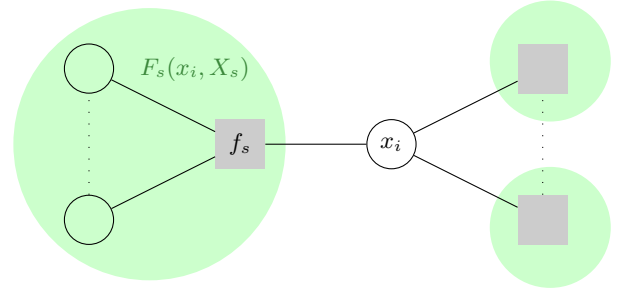


Figure 2: Exemplary factor graph to deduce the messages of the sum-product algorithm.

from all its neighbours in the factor graph:

$$\begin{aligned} p(x_i) &= \sum_{X \setminus x_i} \left[\prod_{s \in ne(x_i)} F_s(x_i, X_s) \right] \\ &= \prod_{s \in ne(x_i)} \left[\sum_{X_s} F_s(x_i, X_s) \right] =: \prod_{s \in ne(x_i)} \mu_{f_s \rightarrow x_i} \end{aligned} \quad (5)$$

The factors $F_s(x_i, X_s)$ represent subtrees in the factor graphs and can thus be factorized again:

$$F_s(x_i, X_s) = f_s(x_i, X_s) G_1(x_1, X_{s_1}) \dots G_M(x_M, X_{s_M}). \quad (6)$$

The factors $G_1(x_1, X_{s_1}) \dots G_M(x_M, X_{s_M})$ represent here the subgraphs of all neighbours of f_s except from x_i and are illustrated by red circles in figure 3.

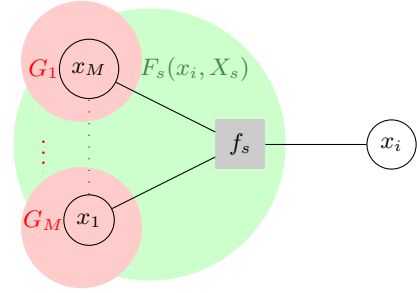


Figure 3: Exemplary factor graph to deduce the messages of the sum-product algorithm.

With that, the messages defined in equation 6 can be further evaluated to

$$\begin{aligned} \mu_{f_s \rightarrow x_i}(x_i) &= \sum_{X_s \setminus x_i} f_s(x_i, X_s) \prod_{m \in ne(f_s) \setminus x_i} \left[\sum_{X_{s_m}} G_m(x_m, X_{s_m}) \right] \\ &= \sum_{\mathbf{x}_s \setminus x_i} f_s(x_i, X_s) \prod_{m \in ne(f_s) \setminus x_i} \mu_{x_m \rightarrow f_s}(x_m). \end{aligned} \quad (7)$$

We have now defined messages from factor to variable nodes. $\mu_{x_m \rightarrow f_s}(x_m) := \sum_{X_{s_m}} G_m(x_m, X_{s_m})$. Going one step further yields an expression for messages from variables to factors as well. Analogously to before, $G_m(x_m, X_{s_m})$ can be factorised again to

$$G_m(x_m, X_{s_m}) = \prod_{l \in ne(x_m) \setminus f_s} F_l(x_m, X_{s_{ml}}), \quad (8)$$

as visualised in figure 4.

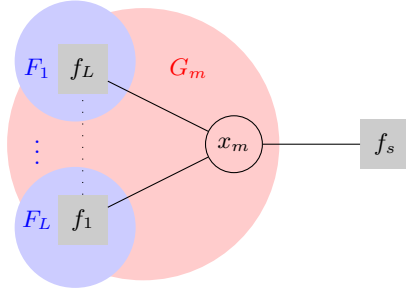


Figure 4: Exemplary factor graph to deduce the messages of the sum-product algorithm.

Inserting this factorisation into the definition of $\mu_{x_m \rightarrow f_s}(x_m)$ and exchanging summation and multiplication reveals an explicit formula for this message:

$$\begin{aligned} \mu_{x_m \rightarrow f_s}(x_m) &= \sum_{X_{s_m}} \left[\prod_{l \in \text{ne}(x_m) \setminus f_s} F_l(x_m, X_{s_m l}) \right] \\ &= \prod_{l \in \text{ne}(x_m) \setminus f_s} \left[\sum_{X_{s_m l}} F_l(x_m, X_{s_m l}) \right] \quad (9) \\ &= \prod_{l \in \text{ne}(x_m) \setminus f_s} \mu_{f_l \rightarrow x_m}(x_m) \end{aligned}$$

We have now induced how a marginal distribution $p(x_i)$ can be expressed by local messages passed between factors and variables through the factor graph. For leave nodes those messages can be stated explicitly. For a variable node as leave of the factor graph equation 9 shows that the message, the node sends to its only neighbor, is given by

$$\mu_{x \rightarrow f}(x) = 1. \quad (10)$$

Similarly, from equation 8 then concludes the initial value

$$\mu_{f \rightarrow x}(x) = f(x) \quad (11)$$

for a message from a leave factor node to its neighbouring variable node.

Having explicit start values for messages from the leaves, messages can be propagated through the factor graphs to an arbitrarily chosen root node. This root node has then received messages from all its neighbours and can thus calculate its marginal distribution according to 5. Analogously, passing all messages back from the root to all leave nodes enables to calculate the marginals of all variables in the factor graph. For this calculation two times the number of edges in the factor graph many messages have to be computed, which is linear in the number of variables of the model. **todo: marginal of a factor formula, connection to EM-algorithm, continuous RVs**

2.3 Maximum a posteriori estimation

Next to computing marginal distributions, a frequent task is to find maximum a posteriori estimate. The aim is to find values for all random variables, which maximize the joint probability distribution and can be formulated mathematically as

$$X^* = \arg \max_X p(X) = \arg \max_X \ln(p(X)). \quad (12)$$

The value for the maximal joint probability is given by

$$p(X^*) = \max_X \ln(p(X)). \quad (13)$$

The actual algorithm works similarly to the sum-product algorithm from the previous chapter, but with the exception that maximizations instead of sums over the random variables are performed. Analogously to before, the idea behind is to insert the factorized expression of the joint distribution from equation 6 into the problem formulation from equation 13 and exchange maximization and multiplication.

The second step of equation 12 contains a transformation, which allows a simplification. Since the logarithm is a strictly monotonic increasing function, it does not change the optimum X^* but transforms the products in the messages to sums. Thus, the algorithm derived in this section is not called "max-product", which the analogy to the max-sum algorithm would suggest, but "max-sum" algorithm.

Applying all discussed modification to the formulas deduced in the last chapter leads to the following messages:

$$\mu_{f_s \rightarrow x_i}(x_i) = \max_{X_s \setminus x_i} \left[f_s(x_i, X_s) \sum_{m \in \text{ne}(f_s) \setminus x_i} \mu_{x_m \rightarrow f_s}(x_m) \right] \quad (14)$$

$$\mu_{x_m \rightarrow f_s}(x_m) = \sum_{l \in \text{ne}(x_m) \setminus f_s} \mu_{f_l \rightarrow x_m}(x_m) \quad (15)$$

Also the message initializations for leave nodes from equation 16 and 17 must be adjusted with the logarithm, resulting in

$$\mu_{x \rightarrow f}(x) = 0, \quad (16)$$

$$\mu_{f \rightarrow x}(x) = \ln(f(x)). \quad (17)$$

By propagation messages from all leave nodes to an arbitrarily chosen root node, the maximal value of the joint probability distribution can be calculated in that root node. The message passing is again done in a fashion, that a node can forward its message to its parent node as soon it has received messages from all children.

The missing step is to find the variable configuration for the maximal joint probability. The central maximization is performed in equation 14 over the variables of each factor before its message is passed further to the root. By saving the maximal configuration of each factor at this step, the values of the variables which give rise to the maximal joint probability can be obtained. In contrast to the sum-product algorithm, the max-sum algorithm does not require to propagate messages back from the root to all leaves.

2.4 Inference in general graphs

loopy belief propagation

3. APPROXIMATE INFERENCE

This chapter presents by the expectation propagation (EP) algorithm a method for approximate inference. In the following section the algorithm and its methodology are explained at first. Secondly it is applied to graphical models which came already up in the previous chapters. At last a general discussion and outlook of EP is given.

Known as expectation propagation, similar to variational inference Minimize KL-divergence using moment matching. Interesting properties different from standard VI

3.1 Methodology

Similar to variational inference (VI) in EP the true distribution $p(\mathbf{z})$ of a model is approximated by a simpler distribution $q(\mathbf{z})$. The approximation is computed and measured by minimizing the Kullback-Leibler divergence of the two functions. However, while in VI $KL(q||p)$ is considered, EP swaps the distributions and aims at minimizing $KL(p||q)$.

For $q(\mathbf{z})$ a member of the exponential family with the form

$$q(\mathbf{z}) = h(\mathbf{z})g(\eta) \exp(\eta^T(u(\mathbf{z}))). \quad (18)$$

With that the KL divergence becomes

$$KL(p||q) = -\ln g(\eta) - \eta^T \mathbf{E}_{p(\mathbf{z})}[u(\mathbf{z})] + \text{const.} \quad (19)$$

The distribution $q(\mathbf{z})$ which minimizes this KL divergence from equation 19 can be constructed by moment matching. If $q(\mathbf{z})$ is a Gaussian for example, its mean corresponds to the mean of $p(\mathbf{z})$ and its covariance is equal to the covariance of $p(\mathbf{z})$.

For the concrete EP algorithm we consider a model of hidden variables and parameters θ and observed data \mathbf{D} . Similarly to the previous chapter we assume that the joint distribution can be written as a product of factors f_i :

$$p(\mathbf{D}, \theta) = \prod_i f_i(\theta). \quad (20)$$

The goal is then to approximate the posterior distribution $p(\theta|\mathbf{D})$, which can be expressed by the joint distribution and the evidence $p(\mathbf{D})$:

$$p(\theta|\mathbf{D}) = \frac{1}{p(\mathbf{D})} \prod_i f_i(\theta), \quad (21)$$

with

$$p(\mathbf{D}) = \int \prod_i f_i(\theta) d\theta. \quad (22)$$

evtl. Footnote to reference.

Adapting to the structure of the true posterior distribution, $q(\theta)$ is chosen as a product of factors \tilde{f}_i from a exponential family, which each correspond to one of the true factors f_i : $\frac{1}{Z}$ represents the normalizing factor for the distribution:

$$q(\theta) = \frac{1}{Z} \prod_i \tilde{f}_i(\theta) \quad (23)$$

The factors \tilde{f}_i are refined iteratively. Starting with a rough initialization of all approximating factors \tilde{f}_i , in each iteration on factor is refined closer to its true factor f_i , while the rest of the factors are fixed. Assuming we want to refine factor j , we first define the product of the fixed factors $q^{\setminus j}$ for that:

$$q^{\setminus j}(\theta) = \prod_{i \neq j} \tilde{f}_i(\theta) \quad (24)$$

Then we evaluate a new approximation of the posterior by matching the moments of $f_j q^{\setminus j}(\theta)$:

$$q^{\text{new}}(\theta) \propto f_j q^{\setminus j}(\theta) \quad (25)$$

Finally, we can extract and store a new approximation of \tilde{f}_j from that: $\frac{1}{Z_j}$ again represents a necessary normalizing factor here:

$$\tilde{f}_j = \frac{1}{Z_j} \frac{q^{\text{new}}(\theta)}{q^{\setminus j}(\theta)} \quad (26)$$

3.2 Expectation propagation in graphical models

Whereas general expectation propagation allows the factors of the joint distribution to depend on all parameters and latent variables of the model (see equation 21), in graphical models the factors of the joint distribution depend on a subset of variables θ_i .

$$p(\theta|\mathbf{D}) = \frac{1}{p(\mathbf{D})} \prod_i f_i(\theta_i), \quad (27)$$

An interesting property gets visible if we use a fully factorized approximation function: Then, the EP algorithm becomes equal the sum-product algorithm. In other words, the sum-product algorithm turns out to be a special case of expectation propagation, if using fully independent variables in the approximation function.

With respect to approximation function of the EP algorithm an interpretation can be formulated: Since the sum-product algorithm is an exact inference algorithm, a more flexible approximation function (as for example one with a few dependency connections in the graph left out) can increase the accuracy.

Another approach to incorporate some extent of generalization to expectation propagation is to refine a set of factors together within one iteration of the algorithm. However, for both methods there are still open research question how to find the best amount of independence for the approximation function of EP as well as the optimal grouping of factors for simultaneous refinement.

3.3 Discussion

A key disadvantage of expectation propagation is that is not guaranteed to converge. For the exponential family the resulting divergence function is proven to have stationary point, but the iteration of EP might not reach it.

Variational inference in contrast provides a lower bound in each optimization step and thus guarantees to converge. However, EP can nevertheless be a powerful inference method. Because variational inference operates with a lower bound it often underestimates the variance of the true distribution. Therefore, if expectation propagation converges, it often outperforms other variational methods. On the other hand, if the algorithm does not converge, it might be for a good reason an the chosen exponential family approximates the true posterior only poorly.

Another aspect leading to conceptually different behavior between expectation propagation and variational inference comes from the approach, the KL divergence is minimized. Expectation propagation minimizes $KL(p||q)$, while variational inference aims at reducing $KL(q||p)$. For multimodal distribution this leads to significantly worse results for expectation propagation. Since the EP algorithm tries to incorporate all modes of the underlying probability distribution it can approximate function with more than one modes only badly. This is visualized in figure 5, which is taken from [1].

4. SUMMARY AND OUTLOOK

4.1 Subsection

blabla with 2 sources [2; 1].

5. REFERENCES

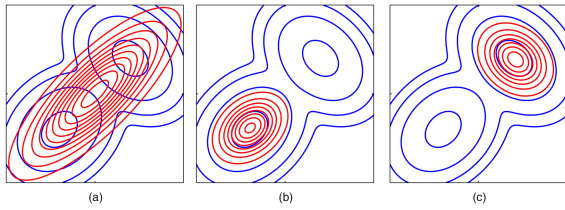


Figure 5: Approximation of a multimodal distribution by a) minimizing $KL(p||q)$, b) minimizing $KL(q||p)$ and c) same approach as b) but with a different result.

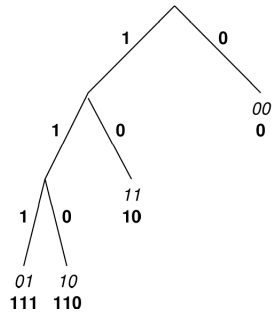


Figure 6: Tree

- [1] C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer-Verlag New York, Inc., 2006.
- [2] K. P. Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012.

Table 1: Example table	
Column 1	Column 2
0	1