

An Auxiliary Variational Method

Felix V. Agakov¹ and David Barber²

¹ University of Edinburgh, 5 Forrest Hill, EH1 2QL Edinburgh, UK
felixa@inf.ed.ac.uk, <http://anc.ed.ac.uk>

² IDIAP, Rue du Simplon 4, CH-1920 Martigny Switzerland,
david.barber@idiap.ch

Abstract. An attractive feature of variational methods used in the context of approximate inference in undirected graphical models is a rigorous lower bound on the normalization constants. Here we explore the idea of using augmented variable spaces to improve on the standard mean-field bounds. Our approach forms a more powerful class of approximations than any structured mean field technique. Moreover, the existing variational mixture models may be seen as computationally expensive special cases of our method. A byproduct of our work is an efficient way to calculate a set of mixture coefficients for any set of tractable distributions that principally improves on a flat combination.

1 Introduction

Probabilistic treatment of uncertainty provides a principled way of reasoning in stochastic domains. Unfortunately, mathematical consistency often comes at a price of inherent intractability of many interesting models, such as Boltzmann machines

$$p(\mathbf{x}) = \exp\{-E(\mathbf{x})\}/Z, \quad Z = \sum_{\mathbf{x}} \exp\{-E(\mathbf{x})\}. \quad (1)$$

In general the complexity of evaluating the partition function Z is exponential in the size of the largest clique in the associated junction tree. For dense models the exact evaluations are in general computationally infeasible, and approximations need to be considered. In this paper we focus on computation of lower bounds on Z , which may also be used to approximate formally intractable marginals.

Variational approximations have been widely used in physics and engineering and more recently applied to graphical modeling (e.g. [5]). In this context they are typically used to obtain rigorous (but relatively simple) bounds on the normalizing constant. A popular class of such methods is based on the KL-divergence

$$KL(q(\mathbf{x})\|p(\mathbf{x})) = \langle \log q(\mathbf{x}) \rangle_{q(\mathbf{x})} - \langle \log p(\mathbf{x}) \rangle_{q(\mathbf{x})} \geq 0, \quad (2)$$

where $\langle \dots \rangle_{q(\mathbf{x})}$ denotes an average over $q(\mathbf{x})$, and the bound is saturated if and only if $q(\mathbf{x}) \equiv p(\mathbf{x})$. In the case of the Boltzmann distribution (1), non-negativity of (2) yields the well-known class of lower bounds

$$\log Z \geq -\langle \log q(\mathbf{x}) \rangle_{q(\mathbf{x})} - \langle E(\mathbf{x}) \rangle_{q(\mathbf{x})}, \quad (3)$$

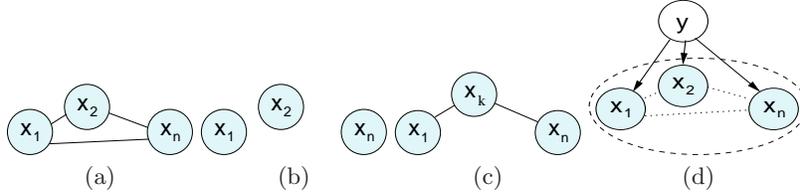


Fig. 1. (a) A fully connected network representing the intractable $p(\mathbf{x})$; (b) standard mean field model $q_{MF}(\mathbf{x})$; (c) structured mean field model $q_{SMF}(\mathbf{x})$; (d) a mixture of mean field models [all the variables \mathbf{x} are coupled through the mixture label y .]

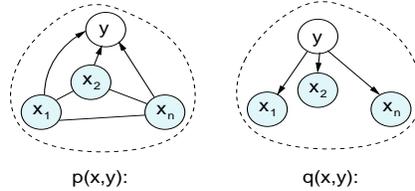


Fig. 2. An auxiliary MF model. The target $p(\mathbf{x}, y)$ is approximated by $q(\mathbf{x}, y)$, which is structured *in the augmented space*. [Note that the marginal $p(\mathbf{x})$ expressed from $p(\mathbf{x}, y)$ is identical to the original fully connected pairwise distribution shown on Figure 1 (a).]

where $q(\mathbf{x})$ is typically restricted to lie in a tractable family and varied to obtain the tightest bound within the family. Coupled with an upper bound on Z , expression (3) may be used for bounding marginals of $p(\mathbf{x})$. This procedure may also be used to optimize a lower bound on the marginal likelihood in partially observable models, which is a natural generalization of the EM algorithm [7].

1.1 Existing Variational Approximations

The tractability of the bound (3) depends on the choice of the approximating distribution $q(\mathbf{x})$, which is in the simplest case given by the factorized *mean field* (MF) model $q_{MF}(\mathbf{x}) = \prod_i q(x_i)$ (see Fig. 1 (a), (b)). Factorized approximations may be simple, but inaccurate when $p(\mathbf{x})$ is strongly coupled; moreover, due to uni-modality they may miss a significant mass contributing to Z . One way to go beyond the factorized assumption for $q(\mathbf{x})$ is to consider a *structured mean field* approximation, which also introduces conditional independencies, but retains some of the structure of $p(\mathbf{x})$. Often it is assumed that $q(\mathbf{x})$ has a sparse graphical representation (e.g. it is a (poly)tree, see Figure 1 (c)), which typically leads to an improvement on the bound at a moderate increase in computational cost. Another approach examined recently [4], [6] uses mixtures of mean field type models (see Figure 1 (d)). This is a powerful extension of factorized approximations, since the resulting $q(\mathbf{x})$ is in general multi-modal and coupled in \mathbf{x} (Fig. 1); however, in this case the bound (3) is itself intractable. Known techniques

[6], [4], [3] handle the intractability by effectively using the Jensen’s bound on top of (3), which is computationally costly and numerically unstable in practice unless *all* the mixture components $q(x|y)$ have the same structure.

2 Auxiliary Variational Method

Intuitively, evaluation of the bound (3) on $\log Z$ for variational mixture approximations requires minimization of the KL-divergence between two fully-connected distributions (see (2)). However, computationally it could be useful to retain a sparser structural form of $q(x, y)$ and use it as an approximation. To do this, we introduce *auxiliary variables* y to the target distribution in such a way that the marginal $\tilde{p}(x)$ of the augmented model $p(x, y)$ has the same graphical structure as the original target $p(x)$ (see Fig. 2). Then we minimize the KL-divergence between $q(x, y)$ and $p(x, y)$ in the joint variable spaces. This case is different from standard structured approximations, as all the variables x of the marginal $\langle q(x|y) \rangle_{q(y)}$ remain fully connected. However, similarly to structured mean field methods, the approximation $q(x, y)$ in the joint space is constrained to be sparse.

Another motivation for this work is the reported success of auxiliary sampling techniques, such as *Hybrid Monte-Carlo* or the *Swendsen-Wang* [8] algorithms. It has been shown that by augmenting the original variable space with auxiliary variables and sampling from joint distributions in the augmented spaces, one can achieve a significant improvement over standard MCMC approaches. The purpose of the auxiliary variables in this context is to capture (structural) information about clusters of correlated variables. It may therefore be hoped that an *auxiliary variational* method performing approximations in the augmented space $\{x, y\}$ may improve on simple approximations in $\{x\}$.

2.1 Optimizing the Auxiliary Variational Bound

Let $p(x, y) = p(x)p(y|x)$ define the joint distribution of the original variables x and auxiliary variables y in the augmented $\{x, y\}$ space. From the divergence $KL(q(x, y)||p(x, y))$ in the joint space it is easy to obtain an expression for the lower bound on the normalizing constant of (1), which is given by

$$\log Z \geq \sum_y q(y) [\langle -E(x) - \log q(x|y) \rangle_{q(x|y)}] + \tilde{I}, \quad \tilde{I} = \sum_x \sum_y q(x, y) \log \frac{p(y|x)}{q(y)} \quad (4)$$

where $p(y|x)$ is an arbitrary *auxiliary conditional* distribution. Clearly, (4) decomposes as a convex sum of the standard lower bounds with approximations $q(x|y)$ and a lower bound $\tilde{I}(x, y)$ on the mutual information. This may be used to improve on a single best (tractable) approximation $q(x|y)$, which is reconstructed by trivially setting $p(y|x) \equiv p(y)$. Note that (4) is tractable as long as $p(y|x)$, $q(y)$, and $q(x|y)$ are constrained to lie in tractable families – there is no need to use further variational relaxations in this case. One tractable choice for the auxiliary mapping $p(y|x)$ is a Gaussian (in this case $q(x|y)$ should also be

parameterized as y is real-valued). Another case leading to exact computations is obtained when each node y_i in $p(y|x)$ has a small number of x -parents, and $q(x, y)$ is a tree. Some other parameterizations may not lead to exact bounds, but may nevertheless result in efficient and practically useful approximations (see [1] for discussions and derivations of the EM algorithms for some of these cases).

2.2 Specific Auxiliary Representations

Here we briefly look at two useful choices of $p(y|x)$ for pairwise Markov networks with $E(x) \stackrel{\text{def}}{=} x^T W x$. As usual, we assume that $q(x, y)$ is tractable, i.e. the only problematic term in (4) is the auxiliary expectation $\langle \log p(y|x) \rangle_{q(x, y)}$.

Parametric Constraints on the Auxiliary Distributions

If the auxiliary space is given by a single multinomial variable $y \in \{1, \dots, M\}$, a natural choice for $p(y|x)$ is to use a *softmax* type representation

$$p(y_k|x) \propto \exp \left\{ f(x^T u^{(k)} + b^{(k)}) \right\}, \quad U = \{u^{(1)}, \dots, u^{(M)}\} \in \mathbb{R}^{|\mathcal{X}| \times M}, \mathbf{b} \in \mathbb{R}^M \quad (5)$$

where $f(x; U, \mathbf{b})$ is some differentiable function and $p(y_k|x)$ is the probability of the auxiliary variable y being in state k . Unfortunately, if the weight vector $u^{(k)}$ is dense, one may need to bound $\langle \log p(y_k|x) \rangle_{q(x, y)}$ (a cheaper alternative is to approximate such terms as $\log p(y_k|x)_{q(x)}$ or use a multivariate factorial representation of the auxiliary variables $p(y|x) = \prod_i p(y_i|x)$ with the *Gaussian field* [2] approximation of $\langle \log p(y_i|x) \rangle_{q(x|y)}$ – see [1] for details). For dense weights such approximations usually do not lead to significant deviations of the objective and are shown to be both accurate and efficient [1]; however, the optimized function is no longer a strict bound on $\log Z$. We now consider a tractable case when $p(y|x)$ has constraints on the parental structure for each auxiliary variable.

Structural Constraints on Auxiliary Distributions

If $\pi_x(y_i)$ and $\pi_y(y_i)$ are x - and y -parents of y_i in the mapping $p(y|x)$, we get

$$\langle \log p(y|x) \rangle_{q(x, y)} = \sum_{i=1, \dots, |y|} \langle \log p(y_i | \pi_x(y_i), \pi_y(y_i)) \rangle_{q(y_i, \pi_x(y_i), \pi_y(y_i))}. \quad (6)$$

The representational complexity of each conditional in (6) in this case is of the order of $s^{|\pi_x(y_i)| + |\pi_y(y_i)|}$, where s is the number of states (for simplicity assumed to be equal for each variable). Since we are free to choose the form of the distribution $p(y|x)$, we can limit its parental structure so that $|\pi_x(y_i)| + |\pi_y(y_i)|$ is small. Clearly, for discrete variables this allows an exact representation of the conditionals. It is also clear that the computational complexity of evaluating (6) is limited by the cost of marginalization of $|\mathcal{X}| + |\mathcal{Y}| - |\pi_x(y_i)| - |\pi_y(y_i)| - 1$ variables from $q(x, y)$, which is tractable as long as $q(x, y)$ is in a tractable family. E.g., in the special case when $q(y, x) = \prod_{l=1}^{|\mathcal{Y}|} q(y_l) \prod_{j=1}^{|\mathcal{X}|} q(x_j | \rho(x_j))$ is a polytree with a small number of y -parents $\rho(x_j)$, the marginalization is exponential in $|\rho(\pi_x(y_i)) \setminus \{\pi_y(y_i) \cup y_i\}|$, which is acceptable if both $p(y|x)$ and $q(x|y)$ are sparse.

3 Relation to Variational Mixture Models

The existing variational mixture approaches may be viewed as a special case of the *unconstrained* auxiliary formulation (where $p(y|x) \equiv q(y|x)$). In this case (4) is intractable even for factorized mixtures with a few components (as $|x|$ is large). This *requires* further factorized relaxations, such as the ones used by [4], [6], [3]. Also, unless all the components $q(x|y)$ have identical structures, the optimization of the existing bounds may become numerically unstable (as it requires computing non-factorized summations of exponentially small terms [6]).

In our approach, we optimize the bound on $\log Z$ subject to constraints on the auxiliary conditional. By first constraining $p(y|x)$ to be tractable and then optimizing the bound (4), we essentially incorporate the constrained Blahut-Arimoto algorithm into the variational inference framework. Arguably, our approach generalizes variational mixture approximations similarly to the way that the variational EM [7] generalizes the standard EM algorithm: by allowing a flexibility in the choice of $p(y|x)$, it improves numerical stability and helps to significantly simplify the computations. Moreover, our method suggests a way to extend variational mixture approaches to structured auxiliary spaces, which may be useful for boosting the effective number of mixture states.

4 Experimental Results

Throughout the simulations, it was assumed that $p(x)$ is a pairwise Markov net with the energy $E(x) = x^T W x + x^T b$ and $x \in \{-1, 1\}^{|x|}$ (see [1] for details).

Reweighting Structured Representations: A by-product of our framework is a simple and fast way to re-weight any set of differently structured (tractable) distributions $q(x|y)$, which principally improves on trivial combinations. For a uniformly chosen $W \in \mathbb{I}^{10 \times 10}$, we generated $K = 10$ spanning trees with the weights $W^{(m)}$ such that $W_{ij}^{(m)} = W_{ij}$ for all i, j and $m = 1, \dots, K$. Then we optimized (4) for $q(y)$ assuming that $q(x|y_k)$ and $p(y_k|x)$ were fixed (see Fig. 3). In this case the bounds were $L_r \approx 8.38$, $L_u \approx 8.87$, $L_b \approx 8.33$, and $L_{av} \approx 9.52$ for the random, uniform, best single, and auxiliary variational weightings respectively. As expected, the auxiliary method leads to the tightest bound.

Structured Auxiliary Mappings: To investigate an influence of a structure of $p(y|x)$ on (4), we assumed a fixed $q(x, y)$ with $x \in \{-1, 1\}^{|x|}$ and $y \in \{-1, 1\}^{|y|}$. The marginals $q(y_i)$ (flip rate) were fixed to be constant for all the nodes y_i . Fig. 3 (c) shows the improvement \tilde{I} (over the convex combination in (4)) as a function of the flip rate for a model with $|x| = 5$, $|y| = 10$. The total number of x - and y -parents of each y_i was constrained to satisfy $|\pi_x(y_i)| + |\pi_y(y_i)| \leq 4$. We observed that some of the optimal auxiliary mappings were close to the theoretically optimal (but generally intractable) $q(x|y)$, though the choice of the structure proved to be important (see [1] for details).

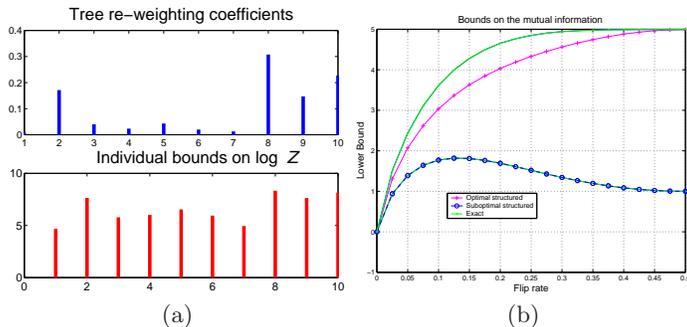


Fig. 3. (a) *Top*: weights of the fixed structured approximations; *Bottom*: bounds on $\log Z$ for each approximation; (b) influence of the structure of $p(y|x)$ on \tilde{I} .

5 Summary

We have presented an approach that generalizes the standard KL- variational procedure to the use of auxiliary variables, which provides a systematic improvement over standard structured approximations. We have also showed that the variational mixture approximations could be seen as special and more computationally expensive cases of our approach. Finally, we showed that our method can be easily generalized to factorial and structural state representations. One way to use it in practice is to find weightings for any set of tractable distributions.

References

1. Agakov, F. V. and Barber, D. (2004). An Auxiliary Variational Method. Technical report, EDI-INF-RR-0205, School of Informatics, University of Edinburgh.
2. Barber, D. and Sollich, P. (2000). Gaussian fields for approximate inference. In *Neural Information Processing Systems 12*. The MIT Press.
3. El-Hay, T. and Friedman, N. (2002). Incorporating Expressive Graphical Models in Variational Approximations: Chain-Graphs and Hidden Variables. In *UAI*.
4. Jaakkola, T. S. and Jordan, M. I. (1998). Improving the Mean Field Approximation via the Use of Mixture Distributions. In Jordan, M. I., editor, *Learning in Graphical Models*. Kluwer Academic Publishers.
5. Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., and Saul, L. K. (1998). An Introduction to Variational Methods for Graphical Models. In Jordan, M. I., editor, *Learning in Graphical Models*, chapter 1. Kluwer Academic Publishers.
6. Lawrence, N. D., Bishop, C. M., and Jordan, M. I. (1998). Mixture Representations for Inference and Learning in Boltzmann Machines. In *UAI: Proceedings of the 14th Conference*.
7. Neal, R. M. and Hinton, G. E. (1998). A View of the EM Algorithm that Justifies Incremental, Sparse, and Other Variants. In Jordan, M. I., editor, *Learning in Graphical Models*, chapter 1. Kluwer Academic Publishers.
8. Swendsen, R. and Wang, J.-S. (1987). Nonuniversal critical dynamics in Monte Carlo simulations. *Physical Review Letters*, 58:86–88.