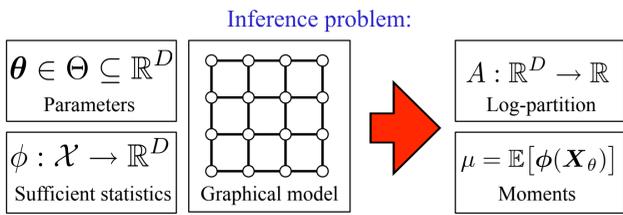


Optimization of Structured Mean Field Objectives

Alexandre Bouchard-Côté* Michael I. Jordan*,†

* Computer Science Division † Department of Statistics
University of California at Berkeley



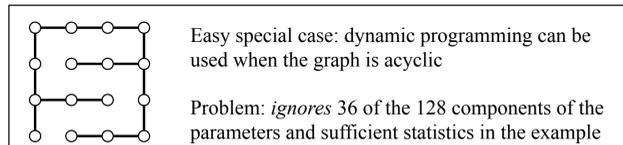
Often intractable (e.g. NP-complete for non-planar Ising models)

$$\mathbb{P}(\mathbf{X}_\theta \in B) = \int_B \exp\{\langle \phi(x), \theta \rangle - A(\theta)\} \nu(dx)$$

$$A(\theta) = \log \int_{\mathcal{X}} \exp\{\langle \phi(x), \theta \rangle\} \nu(dx)$$

$D = \# \text{ vertex} \times \#\{0,1\} + \# \text{ edges} \times \#\{00,01,10,11\}$

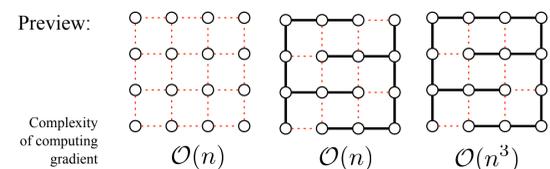
Overview



Structured mean field harnesses an acyclic subgraph, but also takes into account all components

Question: how to choose the acyclic subgraph?

- Adding an edge in the subgraph can only increase quality
- But what is the **impact on computational complexity?**



First result: dichotomy in terms of a graph property, ν -acyclic and b -acyclic subgraphs

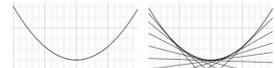
Second result: improved algorithm in the b -acyclic subgraph case

Background

Tool: Legendre-Fenchel transformation

$$f^*(x) = \sup\{\langle x, y \rangle - f(y) : y \in \text{dom}(f)\}$$

Theorem: If f is convex and lower semi-continuous, $f = f^{**}$

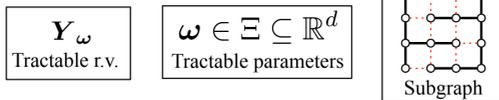


Step 1: Express inference as a constrained optimization problem using convex duality:

$$A(\theta) = \sup\{\langle \theta, \mu \rangle - A^*(\mu) : \mu \in \mathcal{M}\}$$

$$\mathcal{M} = \{\mu \in \mathbb{R}^D : \exists \theta \in \Theta \text{ s.t. } \mathbb{E}[\phi(\mathbf{X}_\theta)] = \mu\}$$

Step 2: Relax the optimization problem using a subset of the initial exponential family (defined by a subgraph)



$$\hat{A}(\theta) = \sup\{\langle \theta, \mu \rangle - A^*(\mu) : \mu \in \mathcal{M}_{\text{MF}}\}$$

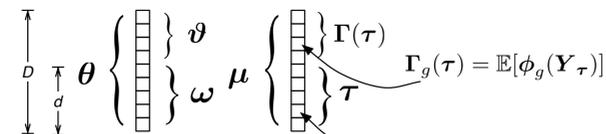
$$\mathcal{M}_{\text{MF}} = \{\mu \in \mathcal{M} : \exists \omega \in \Xi \text{ s.t. } \mathbb{E}[\phi(\mathbf{Y}_\omega)] = \mu\}$$

Consequence: on $\mu \in \mathcal{M}_{\text{MF}}$, $A^*(\mu)$ is tractable

Step 3: Solve the simplified optimization problem

$$\hat{A}(\theta) = \sup\{\langle \omega, \tau \rangle + \langle \vartheta, \Gamma(\tau) \rangle - A_0^*(\tau) : \tau \in \mathcal{N}\}$$

realizable moments in the subgraph



Necessary optimality condition:

$$0 = \omega + J(\tau)\vartheta - \nabla A_0^*(\tau)$$

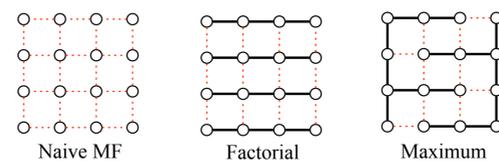
$$\tau = \underbrace{\nabla A_0}_{\text{Easy}}(\omega + \underbrace{J(\tau)\vartheta}_{?}) \quad J = \left(\frac{\partial \Gamma_g}{\partial \tau_f} \right)_{f,g}$$

Dichotomy of tractable mean field subgraphs

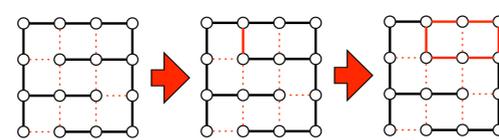
Definition: an acyclic subgraph with edges $E' \subseteq E$ is ...

- ν -acyclic, if for all $e \in E$, $E' \cup \{e\}$ is still acyclic
- b -acyclic, otherwise

Examples of ν -acyclic graphs



Example of a b -acyclic graph



Bag of tricks

Properties of $A(\theta)$

$$\nabla A(\theta) = \mathbb{E}[\phi(\mathbf{X}_\theta)]$$

$$H(A(\theta)) = \text{Var}[\phi(\mathbf{X}_\theta)]$$

$$\text{Var} \geq 0 \implies A(\theta) \text{ is convex}$$

$$\nabla A^* = \nabla A^{-1}$$

when the family is regular and minimal

Hammersley-Clifford Theorem

Consequence: if a, b belong to different cliques,

$$\mathbf{Y}_a \perp\!\!\!\perp \mathbf{Y}_b$$

Chain rule for Jacobian matrices

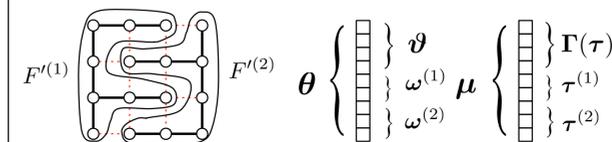
$$I = \left(\frac{\partial Z^{[g]}}{\partial \theta_h^{[g]}} \right)_{g,h}; \quad K = \left(\frac{\partial \theta_h^{[g]}}{\partial \tau_f} \right)_{h,f}$$

$$\implies J = K^T I^T$$

ν -acyclic subgraphs

Connected component decomposition:

$$\tau = \begin{pmatrix} \tau^{(1)} \\ \tau^{(2)} \end{pmatrix} = \begin{pmatrix} \nabla A_0^{(1)}(\omega^{(1)} + J^{(1)}(\tau)\vartheta) \\ \nabla A_0^{(2)}(\omega^{(2)} + J^{(2)}(\tau)\vartheta) \end{pmatrix}$$



Form of J :

$$J_{f,g}^{(1)}(\tau) = \frac{\partial}{\partial \tau_f} \mathbb{E}[\phi_g(\mathbf{Y}_\tau)]$$

$$= \frac{\partial}{\partial \tau_f} \mathbb{P}(Y_a = s, Y_b = t)$$

$$= \frac{\partial}{\partial \tau_f} \mathbb{P}(Y_a = s) \mathbb{P}(Y_b = t)$$

$$= \frac{\partial}{\partial \tau_f} \tau_{a,s} \tau_{b,t} = \tau_{b,t}$$

Relation to block Gibbs sampling

$$\mathbf{X}_t^{(2)} | \mathbf{X}_{t-1} \sim \text{MRF}(\omega^{(2)} + B^{(2)}(\mathbf{X}_{t-1})\vartheta)$$

$J(\tau)$	0.1	0.3		
	0.9		0.8	
		0.2		0.6
		0.5	0.2	0.1

$B(\mathbf{X}_t)$	0	1		
	1		1	
			0	1
		0		0

b -acyclic subgraphs

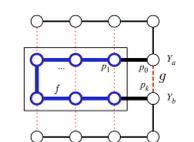
Form of J :

$$J_{f,g}(\tau) = \frac{\partial}{\partial \tau_f} \mathbb{P}(Y_a = s, Y_b = t)$$

$$= \frac{\partial}{\partial \tau_f} \sum_{y_1 \in \mathcal{X}} \dots \sum_{y_{k-1} \in \mathcal{X}} \mathbb{P}(Y_{p_1} = y_1, \dots, Y_{p_{k-1}} = y_{k-1})$$

$$= \frac{\partial}{\partial \tau_f} \mathbb{P}(Y_a = s) \sum_{y_1 \in \mathcal{X}} \mathbb{P}(Y_{p_1} = y_1 | Y_{p_0} = y_0) \sum_{y_2 \in \mathcal{X}} \dots$$

$$= \frac{\partial}{\partial \tau_f} \tau_{a,s} \sum_{y_1 \in \mathcal{X}} \frac{\tau_{(p_0, p_1), (y_0, y_1)}}{\tau_{p_0, y_0}} \sum_{y_2 \in \mathcal{X}} \dots$$



Technique: auxiliary exponential families

For fixed g , construct an exponential families such that its partition function satisfies:

$$Z^{[g]}(\theta^{[g]}) = \sum_{x \in \mathcal{X}^{k-1}} \exp\{\langle \phi(x), \theta^{[g]} \rangle\}$$

$$= \left(\sum_{s'} \tau_{(a, p_1), (s, s')} \right) \sum_{y_1 \in \mathcal{X}} \frac{\tau_{(p_0, p_1), (y_0, y_1)}}{\left(\sum_{s'} \tau_{(p_0, p_1), (y_0, s')} \right)}$$

$$\times \sum_{y_2 \in \mathcal{X}} \dots \sum_{y_{k-1} \in \mathcal{X}} \frac{\tau_{(p_{k-2}, p_{k-1}), (y_{k-2}, y_{k-1})}}{\left(\sum_{s'} \tau_{(p_{k-2}, p_{k-1}), (y_{k-2}, s')} \right)}$$

$$\times \frac{\tau_{(p_{k-1}, p_k), (y_{k-1}, y_k)}}{\left(\sum_{s'} \tau_{(p_{k-1}, p_k), (y_{k-1}, s')} \right)}$$

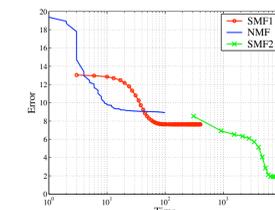
Why? We can get all the derivatives of the log-partition function in one shot using sum-product

$$\frac{\partial Z^{[g]}}{\partial \theta_h^{[g]}} = Z^{[g]} \times \frac{\partial A^{[g]}}{\partial \theta_h^{[g]}} = Z^{[g]} \times \mu_h^{[g]}$$

How can we get the partial derivative with respect to τ ?

Experiments

Adding edges improves the quality of the approximation



Using a b -acyclic subgraph is significantly more expensive

References

D. Barber and W. Wiegand. Tractable variational structures for approximating graphical models. In *Advances in Neural Information Processing Systems*, pages 183–189. Cambridge, MA, 1999. MIT Press.

D. M. Blei and M. I. Jordan. Variational inference for Dirichlet process mixtures. *Bayesian Analysis*, 1:121–144, 2005.

N. De Freitas, P. Højten-Sørensen, M. I. Jordan, and S. Russell. Variational MCMC. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, San Mateo, CA, 2001. Morgan Kaufmann.

D. Geiger, C. Meek, and Y. Wexler. A variational inference procedure allowing internal structure for overlapping clusters and deterministic constraints. *Journal of Artificial Intelligence Research*, 27:1–23, 2006.

A. Globerson and T. Jaakkola. Approximate inference using planar graph decomposition. In *Advances in Neural Information Processing Systems*, Cambridge, MA, 2006. MIT Press.

G. Hua and Y. Wu. Sequential mean field variational analysis of structured deformable shapes. *Computer Vision and Image Understanding*, 101:87–99, 2006.

J. B. Lasserre. Global optimization with polynomials and the problem of moments. *SIAM Journal on Optimization*, 11:796–817, 2000.

T. Minka and Y. Qi. Tree-structured approximations by expectation propagation. In *Advances in Neural Information Processing Systems*, Cambridge, MA, 2003. MIT Press.

C. Peterson and J. R. Anderson. A mean field theory learning algorithm for neural networks. *Complex Systems*, 1: 995–1019, 1987.

L. K. Saul and M. I. Jordan. Exploiting tractable substructures in intractable networks. In *Advances in Neural Information Processing Systems 8*, pages 486–492. Cambridge, MA, 1996. MIT Press.

M. J. Wainwright and M. I. Jordan. Variational inference in graphical models: The view from the marginal polytope. In *Forty-first Annual Allerton Conference on Communication, Control, and Computing*, 2003.

M. J. Wainwright and M. I. Jordan. Log-determinant relaxation for approximate inference in discrete Markov random fields. *IEEE Transactions on Signal Processing*, 54:2099–2109, 2006.

M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1:1–305, 2008.

W. Wiegand. Variational approximations between mean field theory and the junction tree algorithm. In *Proceedings of the Sixteenth Conference on Uncertainty in Artificial Intelligence*, pages 626–633. San Mateo, CA, 2000. Morgan Kaufmann.

J. S. Yedidia, W. T. Freeman, and Y. Weiss. Generalized belief propagation. In *Advances in Neural Information Processing Systems*, pages 689–695. Cambridge, MA, 2001. MIT Press.