

Recursive Cavity Modeling for Estimation of Gaussian MRFs*

Stochastic Systems Group

Jason K. Johnson

October 9, 2002

*/mit/jasonj/Public/SSG-OCT9-02

Overview

- Background
 - Graphical Models (MRFs)
 - Exponential Families
 - Gaussian MRFs
 - Information Geometry and Projections
- Model-Thinning Projections
 - Model Selection by greedy edge-removal procedure.
 - Parameters optimized by Iterative Scaling.
- Recursive Cavity Modeling
 - Nested Dissection
 - Cavity Modeling
 - Blanket Modeling
 - Examples

Graphical Models*

Undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ based upon vertices \mathcal{V} with \mathcal{E} (unordered pairs of vertices).

Random variables $\mathbf{x} = (x_i, i \in \mathcal{V})$ are said to be *Markov* w.r.t \mathcal{G} when

$$p(x_A, x_B | x_S) = p(x_A | x_S) p(x_B | x_S)$$

for all $A, B, S \subset \mathcal{V}$ where S separates A from B .

Hammersley-Clifford, 71.[†] \mathbf{x} is Markov w.r.t. \mathcal{G} if and only if $p(\mathbf{x})$ factors according to \mathcal{G} as

$$p(\mathbf{x}) = \frac{1}{Z(\psi)} \prod_{c \in \mathcal{C}} \psi_c(x_c)$$

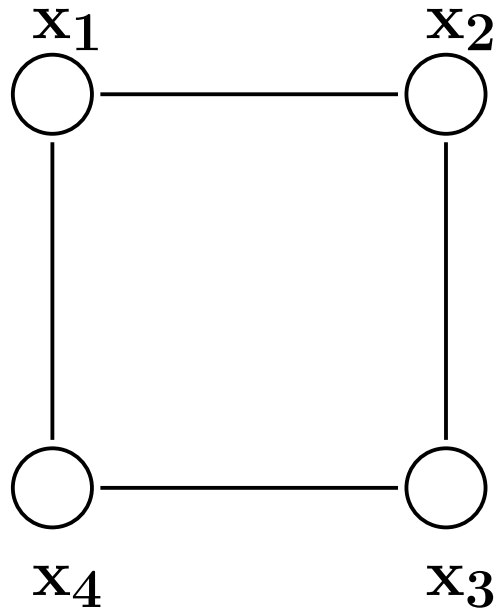
with positive *potential functions* ψ and $Z(\psi)$ is normalization constant.

Markov structure of random process \mathbf{x} allows for compact specification of $p(\mathbf{x})$ as graphical models.

*Lauritzen, 96; Jordan, 99.

[†]Grimmett, 73.

Example MRF



Graph Factorization

$$p(x) \propto \psi_1(x_1)\psi_2(x_2)\psi_3(x_3)\psi_4(x_4) \\ \psi_{1,2}(x_1, x_2)\psi_{2,3}(x_2, x_3) \\ \psi_{3,4}(x_3, x_4)\psi_{4,1}(x_4, x_1)$$

Conditional Independence

$$p(x_{1,3}|x_{2,4}) = p(x_1|x_{2,4})p(x_3|x_{2,4})$$

$$p(x_{2,4}|x_{1,3}) = p(x_2|x_{1,3})p(x_4|x_{1,3})$$

Exponential Families*

Specified by a *base measure* $q(x) > 0$ and a set of *sufficient statistics* $t(x)$ both defined over some specified state-space X . We take $X = \mathbf{R}^n$ so that model is specified by pdf of the form

$$f(x; \theta) = q(x) \exp\{\theta \cdot t(x) - \varphi(\theta)\}$$

where the *cumulant function* $\varphi(\theta)$ is the normalization constant

$$\varphi(\theta) = \log \int q(x) \exp\{\theta \cdot t(x)\} dx$$

Only consider *admissible* parameters Θ s.t. pdf is normalizable $\varphi(\theta) < \infty$. The family is *regular* if Θ has non-empty interior. The statistics are *minimal* if the $t(x)$ are linearly-independent. Then, dual parameterization provided by *moment coordinates* $\eta = E_{\theta}\{t(x)\}$ over the set of achievable moments $\eta(\Theta)$.

*Chentsov, 66; Barndorff-Nielsen, 78.

Gaussian Markov Random Fields

Consider Gaussian process $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with mean vector $\boldsymbol{\mu} = E\{\mathbf{x}\}$ and covariance matrix $\boldsymbol{\Sigma} = E\{\mathbf{x}\mathbf{x}'\} - \boldsymbol{\mu}\boldsymbol{\mu}'$.

Information Filter Form. Say that $\mathbf{x} \sim \mathcal{N}^{-1}(\mathbf{h}, \mathbf{J})$ if

$$\begin{aligned}\mathbf{h} &= \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} \\ \mathbf{J} &= \boldsymbol{\Sigma}^{-1}\end{aligned}$$

s.t. density function is parameterized as

$$p(\mathbf{x}) = \exp\left\{-\frac{1}{2}\mathbf{x}'\mathbf{J}\mathbf{x} + \mathbf{h}'\mathbf{x} - \varphi(\mathbf{h}, \mathbf{J})\right\}$$

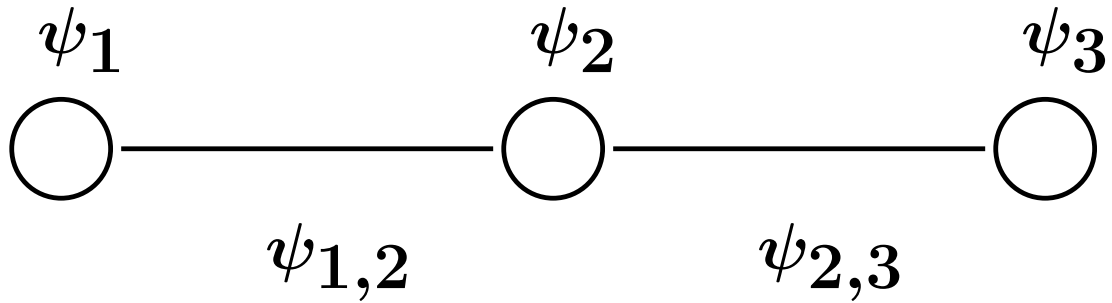
where

$$\varphi(\mathbf{h}, \mathbf{J}) = \frac{1}{2}\{\mathbf{h}'\mathbf{J}^{-1}\mathbf{h} - \log |\mathbf{J}| + n \log 2\pi\}.$$

This is an exponential family model with

$$\begin{aligned}\boldsymbol{\theta} &= (\mathbf{h}, -\mathbf{J}/2) \\ \mathbf{t}(\mathbf{x}) &= (\mathbf{x}, \mathbf{x}\mathbf{x}') \\ \boldsymbol{\eta} &= (\boldsymbol{\mu}, \boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}') \\ \varphi(\boldsymbol{\theta}) &= \varphi(\mathbf{h}, \mathbf{J})\end{aligned}$$

Example GMRF



$$p(x) \propto \psi_1(x_1)\psi_2(x_2)\psi_3(x_3)\psi_{1,2}(x_1, x_2)\psi_{2,3}(x_2, x_3)$$

$$\psi_1(x_1) = \exp\left\{-\frac{1}{2}x_1'J_{1,1}x_1 + h_1'x_1\right\}$$

$$\psi_2(x_2) = \exp\left\{-\frac{1}{2}x_2'J_{2,2}x_2 + h_2'x_2\right\}$$

$$\psi_3(x_3) = \exp\left\{-\frac{1}{2}x_3'J_{3,3}x_3 + h_3'x_3\right\}$$

$$\psi_{1,2}(x_1, x_2) = \exp\{-x_1'J_{1,2}x_2\}$$

$$\psi_{2,3}(x_2, x_3) = \exp\{-x_2'J_{2,3}x_3\}$$

$$h = \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix}, \quad J = \begin{pmatrix} J_{1,1} & J_{1,2} & 0 \\ J_{1,2}' & J_{2,2} & J_{2,3} \\ 0 & J_{2,3}' & J_{3,3} \end{pmatrix}$$

Information Geometry*

Based upon the *Kullback-Leibler divergence*[†], a measure of contrast between probability distributions.

$$D(p\|q) = E_p \left\{ \log \frac{p(\mathbf{x})}{q(\mathbf{x})} \right\}$$

Bregman distance in θ based upon $\varphi(\theta)$,

$$D(\theta^* \|\theta) = \varphi(\theta) - \nabla \varphi(\theta^*) \cdot (\theta - \theta^*)$$

Legendre transform $\varphi^*(\eta)$ of $\varphi(\theta)$:

$$\varphi^*(\eta) = \theta(\eta) \cdot \eta - \varphi(\theta)$$

“Slope transform”

$$\eta(\theta) = \frac{\partial \varphi(\theta)}{\partial \theta}$$
$$\theta(\eta) = \frac{\partial \varphi^*(\eta)}{\partial \eta}$$

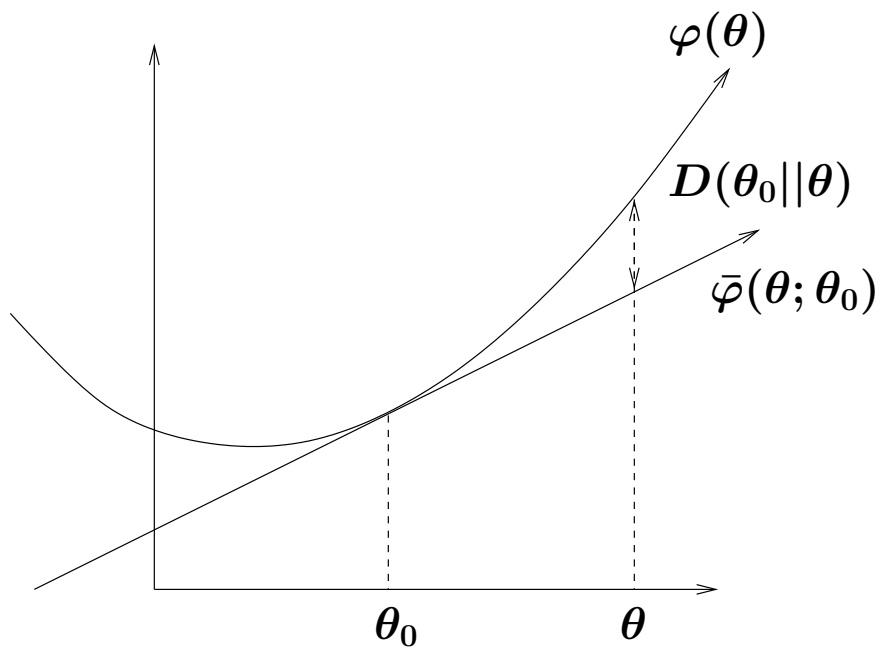
Convex bifunction in $(\eta(p), \theta(q))$,

$$D(\eta \|\theta) = \varphi^*(\eta) + \varphi(\theta) - \eta \cdot \theta$$

*Chentsov, 72; Csiszár, 75; Efron, 78; Amari, 01.

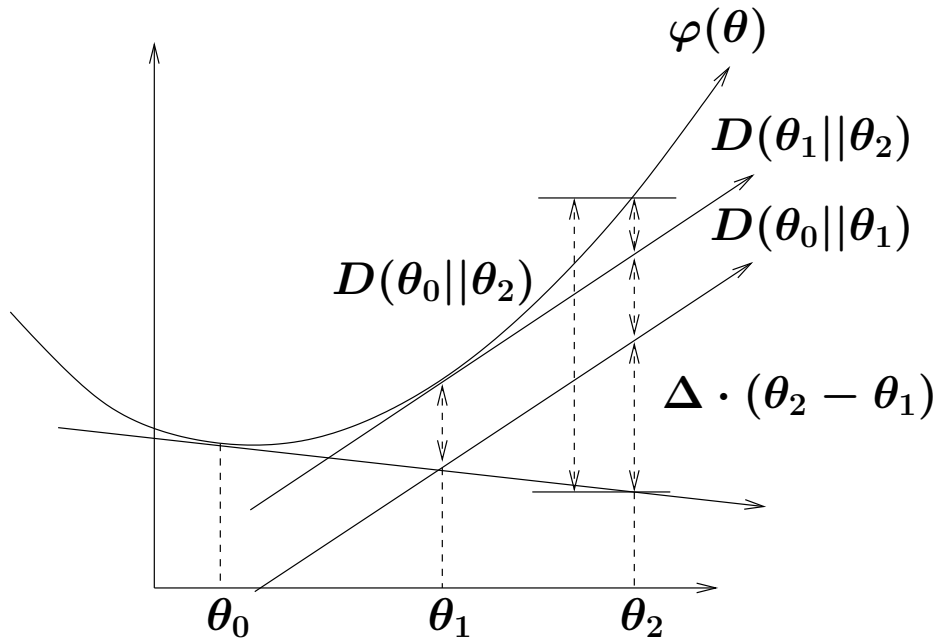
†Kullback and Leibler, 51.

Bregman distance*



*Bregman, 67.

Triangle Relation



$$D(\theta_0||\theta_2) = D(\theta_0||\theta_1) + D(\theta_1||\theta_2) + (\eta_1 - \eta_0) \cdot (\theta_2 - \theta_1)$$

Information Projections

Let \mathcal{F} be a regular exponential family with minimal statistics $t(x)$, exponential coordinates Θ , and moment coordinates $\eta(\Theta)$.

M-projection. Let $p \in \mathcal{F}$, $\mathcal{H} \subset \mathcal{F}$ e-flat submanifold. Exists unique $q^* \in \mathcal{H}$ satisfying the following equivalent conditions:

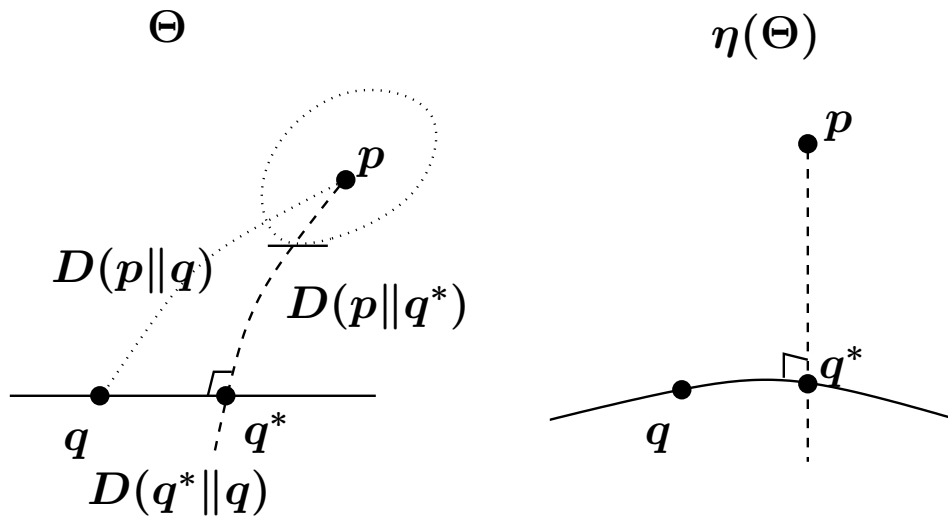
$$(i) \quad D(p\|q^*) = \inf_{q \in \mathcal{H}} D(p\|q)$$

$$(ii) \quad \forall q \in \mathcal{H} : (\eta(p) - \eta(q^*)) \cdot (\theta(q) - \theta(q^*)) = 0$$

$$(iii) \quad \forall q \in \mathcal{H} : D(p\|q) = D(p\|q^*) + D(q^*\|q)$$

We call $q^* = \arg \min_{q \in \mathcal{H}} D(p\|q)$ the *m-projection* of p to \mathcal{H} .

M-projection



$$\frac{\partial}{\partial \theta(q)} D(p||q) = \eta(q) - \eta(p)$$

Dual E-projection

E-projection. Let $q \in \mathcal{F}$, $\mathcal{H}' \subset \mathcal{F}$ m-flat submanifold. Exists unique $p^* \in \mathcal{H}'$ satisfying the following equivalent conditions:

$$(i) \quad D(p^*||q) = \inf_{p \in \mathcal{H}'} D(p||q)$$

$$(ii) \quad \forall p \in \mathcal{H}' : (\eta(p) - \eta(p^*)) \cdot (\theta(q) - \theta(p^*)) = 0$$

$$(iii) \quad \forall p \in \mathcal{H}' : D(p||q) = D(p||p^*) + D(p^*||q)$$

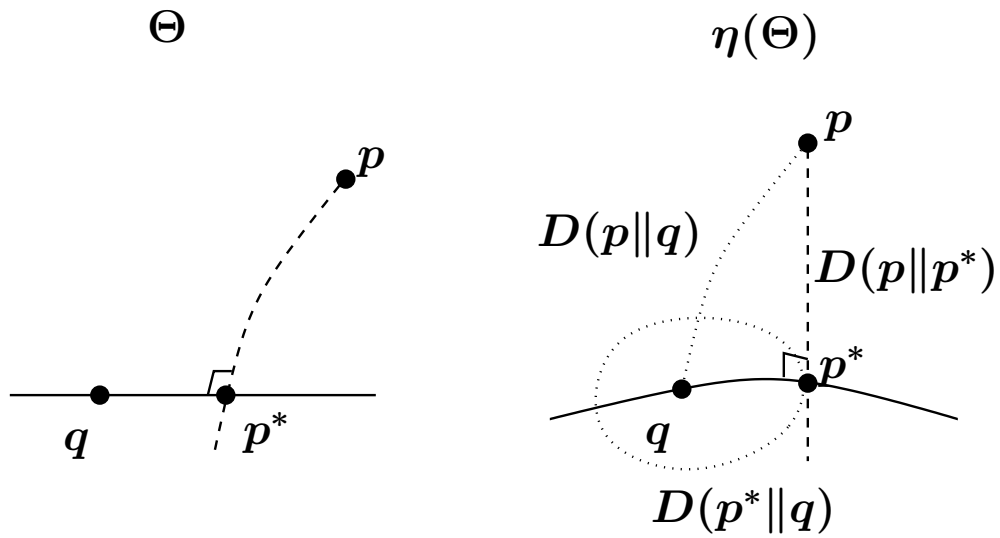
We call $p^* = \arg \min_{p \in \mathcal{H}'} D(p||q)$ the *e-projection* of q to \mathcal{H}' .

Duality. Let \mathcal{H} and \mathcal{H}' be *I-orthogonal* submanifolds such that exists r in intersection and

$$\forall p \in \mathcal{H}', q \in \mathcal{H} : (\eta(p) - \eta(r)) \cdot (\theta(q) - \theta(r)) = 0$$

Then, r is both the m-projection of $p \in \mathcal{H}'$ to \mathcal{H} and the e-projection of $q \in \mathcal{H}$ to \mathcal{H}' .

E-projection



$$\frac{\partial}{\partial \eta(p)} D(p||q) = \theta(p) - \theta(q)$$

Model Thinning

Let $t(x) = (t_{\mathcal{H}}(x), t'_{\mathcal{H}}(x))$, $\theta = (\theta_{\mathcal{H}}, \theta'_{\mathcal{H}})$ and $\eta = (\eta_{\mathcal{H}}, \eta'_{\mathcal{H}})$.

Objective. M-project $p \in \mathcal{F}$ to lower-order exponential family,

$$\mathcal{H} = \{q \in \mathcal{F} \mid \theta'_{\mathcal{H}}(q) = 0\}$$

Dual Problem. E-projection $q \in \mathcal{H}$ to the m-flat submanifold:

$$\mathcal{H}'(p) = \{r \in \mathcal{F} \mid \eta_{\mathcal{H}}(r) = \eta_{\mathcal{H}}(p)\} \quad (1)$$

The latter e-projection problem may be solved by iterative scaling techniques which adjust parameters $\theta_{\mathcal{H}}(q)$ until $\eta_{\mathcal{H}}(q) = \eta_{\mathcal{H}}(p)$ (moment matching).

For GMRF $\mathbf{x} \sim \mathcal{N}^{-1}(h, J)$, impose sparsity on J . Moment-matching gives classical *covariance selection problem* (Dempster, 72).

Iterative Scaling

Alternating e-projections to set of m-flat submanifolds converges to e-projection to intersection (Csiszár, 75). Special case of method of alternating Bregman projections (Bregman, 67).

*Iterative Proportional Fitting.** m-flat submanifolds impose marginal moment constraints specifying marginal distribution $p^*(x_C)$.

$$\psi(x_C) \leftarrow \psi(x_C) \times \frac{p^*(x_C)}{p(x_C)}$$

Covariance Selection.† Updates exponential parameters (h_C, J_C) to impose moment constraints (μ_C^*, Σ_C^*) .

$$\begin{aligned} J_C &\leftarrow J_C + (J_C^* - \hat{J}_C) \\ h_C &\leftarrow h_C + (h_C^* - \hat{h}_C) \end{aligned}$$

where $(h_C^*, J_C^*) = ((\Sigma_C^*)^{-1} \mu_C^*, (\Sigma_C^*)^{-1})$ and $(\hat{h}_C^*, \hat{J}_C^*) = (\Sigma_C^{-1} \mu_C, \Sigma_C^{-1})$ (marginal information models).

*Ireland and Kullback, 68.

†Speed and Kivveri, 86.

Greedy Edge-Removal

Prunes edges from graphical model by forcing selected off-diagonal entries of \mathbf{J} to zero (m-projections implemented by iterative scaling techniques).

Selects weak interactions to prune according to conditional mutual information

$$I(\mathbf{x}_i; \mathbf{x}_j | \mathbf{x}_{\setminus ij}) = -\frac{1}{2} \log \left(1 - \frac{\det \mathbf{J}_{i,j}}{\sqrt{\det \mathbf{J}_{i,i} \det \mathbf{J}_{j,j}}} \right)$$

which gives tractable lower-bound estimate of KL under m-projection.

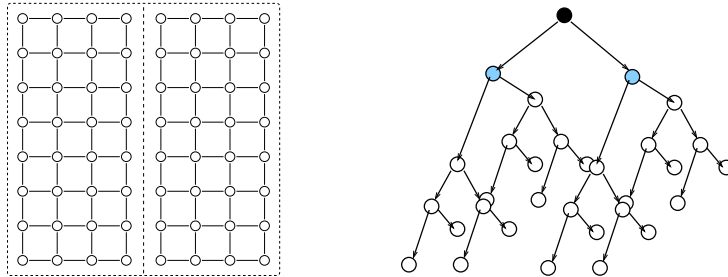
Selects batch $\mathcal{K} \subset \mathcal{V}$ of weakest edges to prune satisfying

$$\sum_{\mathcal{K}} I_{i;j} < \frac{\delta}{|\mathcal{K}|}$$

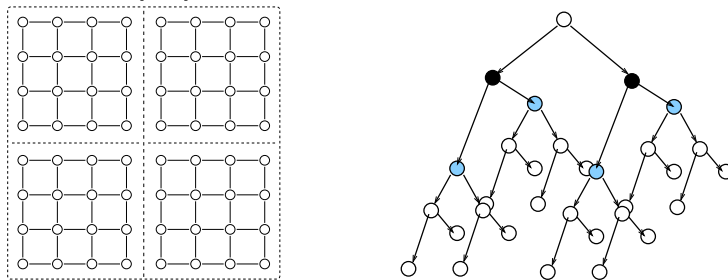
Continues thinning until no more weak interactions relative to δ . Related to Akaike information criterion (Akaike, 74).

Nested Dissection

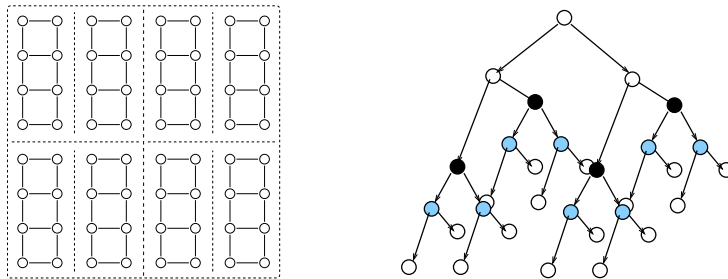
(1) vertical cut.



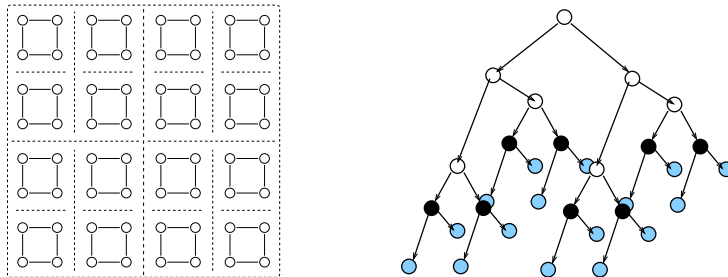
(2) horizontal cut.



(3) vertical cut



(4) horizontal cut.



Variable Elimination

Integrate over subset $\Lambda \subset \mathcal{V}$ of random variables:

$$p(\mathbf{x}_{\setminus\Lambda}) = \int p(\mathbf{x}) d\mathbf{x}_{\Lambda}$$

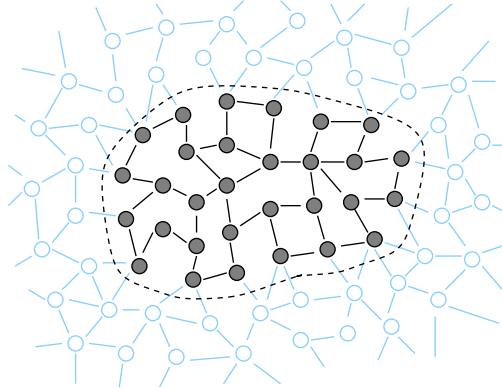
Local parameter update in (\mathbf{h}, \mathbf{J}) representation:

$$\begin{aligned} \mathbf{h}_{\partial\Lambda} &\leftarrow \mathbf{h}_{\partial\Lambda} - \mathbf{J}_{\partial\Lambda,\Lambda} \mathbf{J}_{\Lambda,\Lambda}^{-1} \mathbf{h}_{\Lambda} \\ \mathbf{J}_{\partial\Lambda} &\leftarrow \mathbf{J}_{\partial\Lambda} - \mathbf{J}_{\partial\Lambda,\Lambda} \mathbf{J}_{\Lambda,\Lambda}^{-1} \mathbf{J}_{\Lambda,\partial\Lambda} \end{aligned}$$

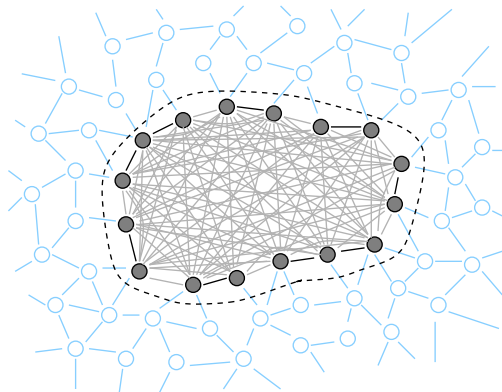
Eliminates vertices in graphical model but adds “fill” edges between neighbors. Only updates local parameters and structure of “boundary” $\partial\Lambda$ of subfield.

Cavity Models (Initialization)

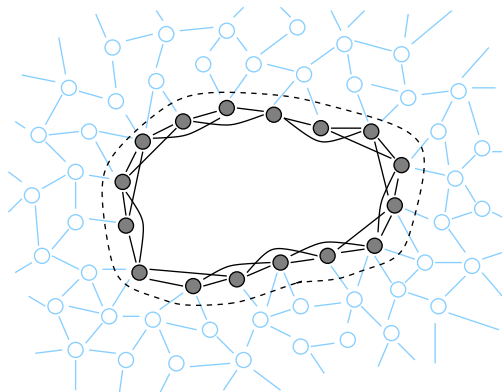
(1) Partial model of subfield (zero boundary).



(2) Elimination gives model of surface.

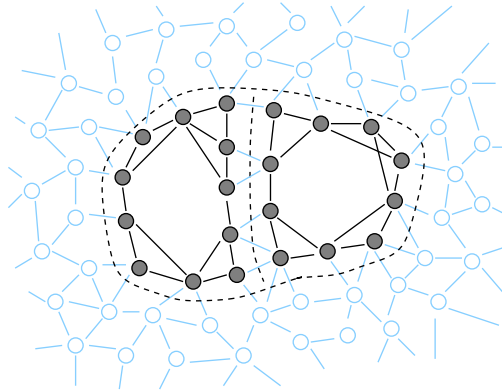


(3) Model thinning gives "cavity model".

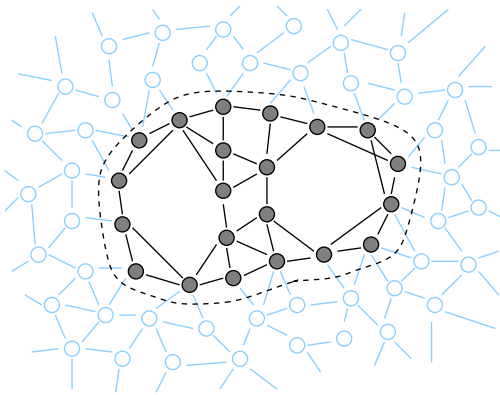


“Upwards” Cavity Modeling

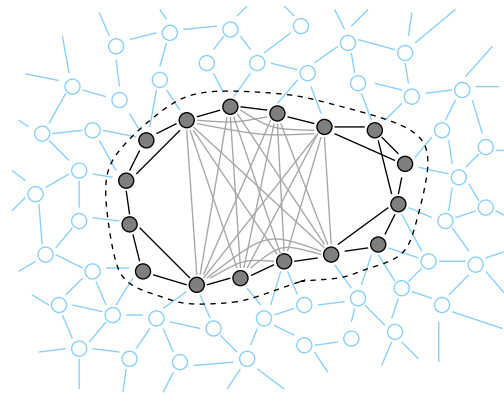
(1) Initialization.



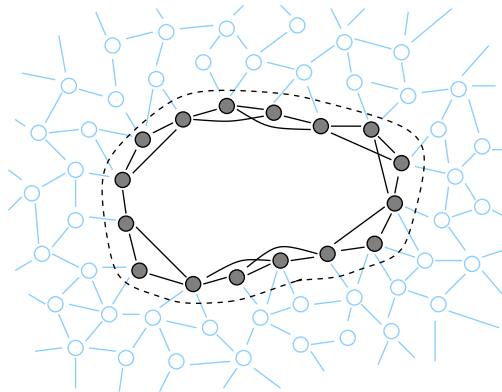
(2) Merge.



(3) Eliminate.

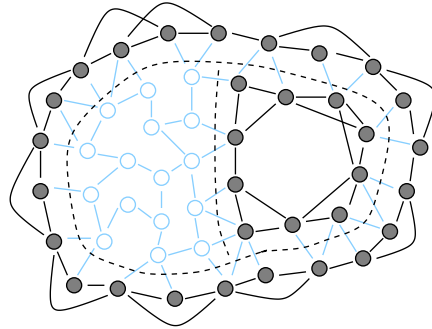


(4) Thin.

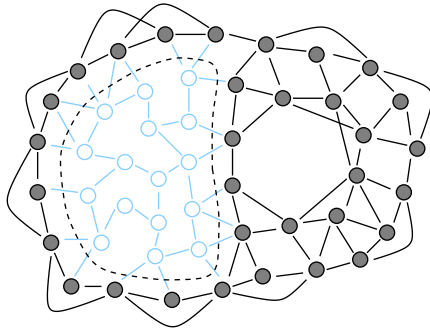


“Downwards” Blanket Modeling

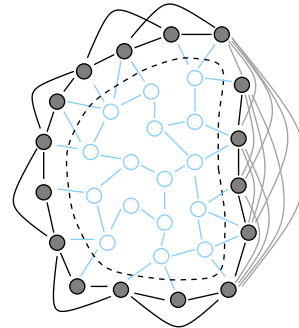
(1) Initialization.



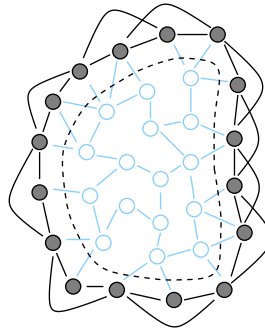
(2) Merge.



(3) Eliminate.



(4) Thin.



Conclusion

RCM appears to provide a powerful and flexible framework for tractable yet near-optimal computation in MRFs.

Much work remains to better characterize performance and explore promising extensions:

- Develop information geometry of RCM.
- Consider more general families of graphical models.
- Employ alternative modeling techniques.
- Applications
 - Model Identification
 - Image Processing
 - Data Compression and Coding
 - Monte-Carlo Simulation

References

- Akaike, 74. A new look at the statistical model identification. *IEEE Trans. Auto. Control*, AC-19:716:723.
- Amari, 01. Information geometry of hierarchy of probability distributions. *IEEE Trans. Inf. Theory*, 47(5):1701-1711.
- Chentsov, 66. A systematic theory of exponential families. *Theory of Prob. and Appl.*, 11.
- Chentsov, 72. Statistical decision rules and optimal inference. *AMS Trans. Math. Mono.*, v.53 (reprint 82).
- Barndorff-Nielsen, 78. *Information and Exponential Families*. John Wiley.
- Bregman, 67. The relaxation method of finding the common point of convex sets. *USSR Comp. Math. and Physics*, 7:200-217.
- Csiszár, 75. I-divergence geometry of probability distributions and minimization problems. *Annals of Prob.*, 3(1):146-158.
- Dempster, 72. Covariance Selection. *Biometrics*, 28(1):157-175.
- Efron, 78. The geometry of exponential families. *Annals of Stat.*, 6(2):362-376.
- Grimmett, 73. A theorem about random fields. *Bull. of London Math. Soc.*, 5:81-84.

Ireland and Kullback, 68. Contingency tables with given marginals. *Biometrika*, 55:179-188.

Jordan (editor), 99. *Learning in Graphical Models*. MIT Press.

Kullback and Leibler, 51. On information and sufficiency. *Annals of Math. Stat.*, 22(1):79-86.

Lauritzen, 96. *Graphical Models*. Oxford University Press.

Speed and Kiiveri, 86. Gaussian Markov distributions over finite graphs. *Annals of Stat.*, 14(1):138-150.