

# Probabilistic Graphical Models

## Lecture 9: Variational Inference Relaxations

Volkan Cevher, Matthias Seeger  
Ecole Polytechnique Fédérale de Lausanne

24/10/2011



# Announcement

- No assignment this week
- Deadline programming assignment: June 18 (next lecture)  
`bayesml09lecture@googlemail.com`

- 1 Structured Mean Field (Variational Bayes)
- 2 Moment Parameters. Variational Relaxations

# Variational Mean Field

$$\log Z \geq \sup_{Q \in \mathcal{Q}} \{E_Q[\Psi(\mathbf{x})] + H[Q(\mathbf{x})]\}$$

- $\mathcal{Q}$ : Tractable subset of all distributions (factorization constraints)

$$\mathcal{Q} = \left\{ Q(\mathbf{x}) = \prod_k Q_k(\mathbf{x}_{S_k}) \right\}, \quad S_k \text{ disjoint}$$

**Tractable?** For any  $k$ ,

$\mathcal{N}_k$ : Factor nodes  $j$  connected to any  $i \in S_k$  ( $S_k \cap C_j \neq \emptyset$ )

$$Q'_k(\mathbf{x}_{S_k}) \propto \exp \left( \sum_{j \in \mathcal{N}_k} E_{Q(\mathbf{x}_{C_j \setminus S_k})} [\Psi_j(\mathbf{x}_{C_j})] \right)$$

tractable to handle

F1

# Variational Mean Field

$$\log Z \geq \sup_{Q \in \mathcal{Q}} \{E_Q[\Psi(\mathbf{x})] + H[Q(\mathbf{x})]\}$$

- $\mathcal{Q}$ : Tractable subset of all distributions (factorization constraints)

$$\mathcal{Q} = \left\{ Q(\mathbf{x}) = \prod_k Q_k(\mathbf{x}_{S_k}) \right\}, \quad S_k \text{ disjoint}$$

**Tractable?** For any  $k$ ,

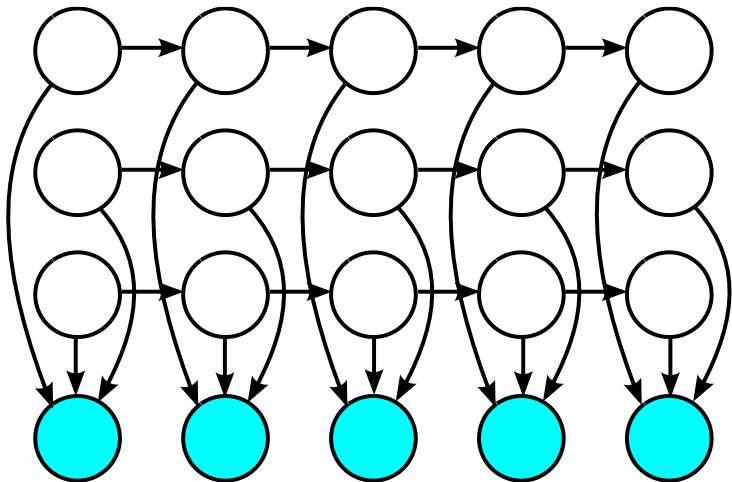
$\mathcal{N}_k$ : Factor nodes  $j$  connected to any  $i \in S_k$  ( $S_k \cap C_j \neq \emptyset$ )

$$Q'_k(\mathbf{x}_{S_k}) \propto \exp \left( \sum_{j \in \mathcal{N}_k} E_{Q(\mathbf{x}_{C_j \setminus S_k})} [\Psi_j(\mathbf{x}_{C_j})] \right)$$

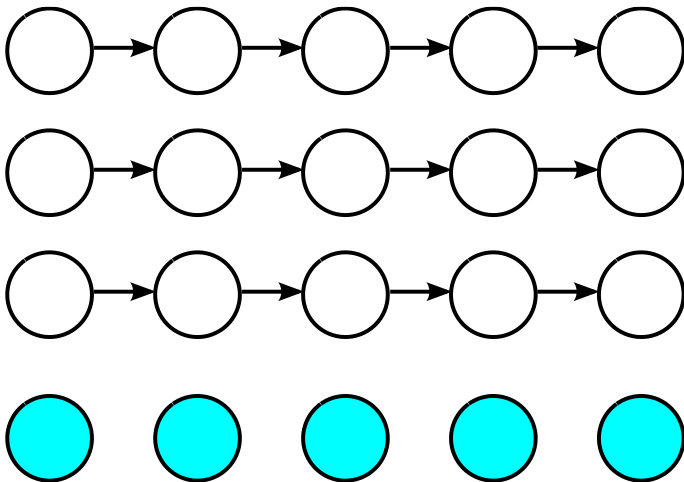
tractable to handle

- $Q(\mathbf{x})$  completely factorized? Naive mean field  
Anything more elaborate? **Structured** mean field

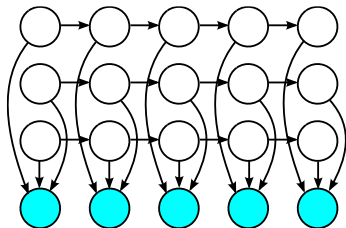
# Factorial Hidden Markov Model



# Factorial Hidden Markov Model



# Factorial Hidden Markov Model

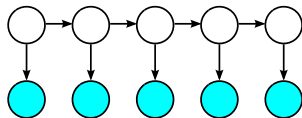
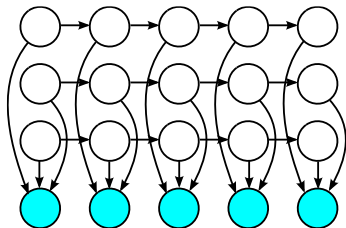


- $S_1 =$  uppermost chain. Update?

F3



# Factorial Hidden Markov Model



- $S_1 =$  uppermost chain. Update?
- $Q(\mathbf{x}_{S_1})$ : Markov chain (variable single node potentials)
  - Double node (transition) potentials of  $Q(\mathbf{x}_{S_k})$ ? Fixed up front!
  - Forward-backward for single node marginals to update  $Q(\mathbf{x}_{S_1})$ . Implementation reduces to single HMM code, called with changing evidence potentials
- **Not** magic, but **as expected**:  
If this does not happen, you made a mistake

# Variational Bayes

- Another instance of re-naming game:  
Nothing else than structured mean field
- Often applied to  $P(\mathbf{x}, \theta | \mathbf{y})$   
( $\mathbf{y}$  observed,  $\mathbf{x}$  latent nuisance,  $\theta$  latent parameters)

# Variational Bayes

- Another instance of re-naming game:  
Nothing else than structured mean field
- Often applied to  $P(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})$   
( $\mathbf{y}$  observed,  $\mathbf{x}$  latent nuisance,  $\boldsymbol{\theta}$  latent parameters)

Expectation maximization

$$\max_{\boldsymbol{\theta}} \log \int P(\mathbf{y}, \mathbf{x} | \boldsymbol{\theta}) d\mathbf{x}$$

$$\geq \max_{\boldsymbol{\theta}, Q(\mathbf{x})} \left\{ \mathbb{E}_Q[\log P(\mathbf{y}, \mathbf{x} | \boldsymbol{\theta})] \right. \\ \left. + H[Q(\mathbf{x})] \right\}$$

Variational Bayes

$$\log \int P(\mathbf{y}, \mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} d\boldsymbol{\theta}$$

$$\geq \max_{Q(\boldsymbol{\theta}), Q(\mathbf{x})} \left\{ \mathbb{E}_Q[\log P(\mathbf{y}, \mathbf{x} | \boldsymbol{\theta})] \right. \\ \left. + H[Q(\mathbf{x})] + H[Q(\boldsymbol{\theta})] \right\}$$

Factorization assumption:  $Q(\mathbf{x}, \boldsymbol{\theta}) = Q(\mathbf{x})Q(\boldsymbol{\theta})$

# Variational Bayes

- Another instance of re-naming game:  
Nothing else than structured mean field

- Often applied to  $P(\mathbf{x}, \theta | \mathbf{y})$   
( $\mathbf{y}$  observed,  $\mathbf{x}$  latent nuisance,  $\theta$  latent parameters)

Expectation maximization

$$\max_{\theta} \log \int P(\mathbf{y}, \mathbf{x} | \theta) d\mathbf{x}$$

$$\geq \max_{\theta, Q(\mathbf{x})} \left\{ \mathbb{E}_Q[\log P(\mathbf{y}, \mathbf{x} | \theta)] \right. \\ \left. + H[Q(\mathbf{x})] \right\}$$

Variational Bayes

$$\log \int P(\mathbf{y}, \mathbf{x} | \theta) d\mathbf{x} d\theta$$

$$\geq \max_{Q(\theta), Q(\mathbf{x})} \left\{ \mathbb{E}_Q[\log P(\mathbf{y}, \mathbf{x} | \theta)] \right. \\ \left. + H[Q(\mathbf{x})] + H[Q(\theta)] \right\}$$

Factorization assumption:  $Q(\mathbf{x}, \theta) = Q(\mathbf{x})Q(\theta)$

- Easy to write generic code (bit like MCMC Gibbs sampling)
- Good approximation?

Can do better today for almost any well-studied model

# Moment Parameters

$$\log Z = \sup_Q \{E_Q[\Psi(\mathbf{x})] + H[Q(\mathbf{x})]\}$$

# Moment Parameters

$$\log Z \geq \sup_{Q \in \mathcal{Q}} \{E_Q[\Psi(\mathbf{x})] + H[Q(\mathbf{x})]\}$$

- $\mathcal{Q}$ : Tractable subset of all distributions (factorization constraints)  
⇒ Seems whole story. What else could there be?

# Moment Parameters

$$\log Z \geq \sup_{Q \in \mathcal{Q}} \{E_Q[\Psi(\mathbf{x})] + H[Q(\mathbf{x})]\}$$

- $\mathcal{Q}$ : Tractable subset of all distributions (factorization constraints)  
 $\Rightarrow$  Seems whole story. What else could there be?
- Consider log-linear models:  $\Psi_j(\mathbf{x}_{C_j}) = \theta_j^T \mathbf{f}_j(\mathbf{x}_{C_j})$ ,  $\theta = (\theta_j)$  F5

$$E_Q[\Psi(\mathbf{x})] = \sum_j \theta_j^T \mu_j, \quad \mu_j := E_Q[\mathbf{f}_j(\mathbf{x}_{C_j})], \quad \mu = (\mu_j)$$

# Moment Parameters

$$\log Z \geq \sup_{Q \in \mathcal{Q}} \{E_Q[\Psi(\mathbf{x})] + H[Q(\mathbf{x})]\}$$

- $\mathcal{Q}$ : Tractable subset of all distributions (factorization constraints)  
 $\Rightarrow$  Seems whole story. What else could there be?
- Consider log-linear models:  $\Psi_j(\mathbf{x}_{C_j}) = \theta_j^T \mathbf{f}_j(\mathbf{x}_{C_j})$ ,  $\theta = (\theta_j)$

$$E_Q[\Psi(\mathbf{x})] = \sum_j \theta_j^T \mu_j, \quad \mu_j := E_Q[\mathbf{f}_j(\mathbf{x}_{C_j})], \quad \mu = (\mu_j)$$

- **Moment parameters**: Under mild assumptions on  $\mathbf{f}_j(\mathbf{x}_{C_j})$ :  
 Just another way (instead of  $\theta$ ) of parameterizing  $P(\mathbf{x})$



# Moment Parameters: Examples

$$P(\mathbf{x}; \boldsymbol{\theta}) = Z^{-1} e^{\boldsymbol{\theta}^T \mathbf{f}(\mathbf{x})}$$

$\boldsymbol{\theta}$	Natural parameters
$\mathbf{f}(\mathbf{x})$	Statistics, representation
$\boldsymbol{\mu} = E_{\boldsymbol{\theta}}[\mathbf{f}(\mathbf{x})]$	Moment parameters

Representation **minimal**: For every  $\mathbf{z} \neq \mathbf{0}$ , there is  $\mathbf{x}$ :

$$\mathbf{z}^T (\mathbf{f}(\mathbf{x})^T \mathbf{1})^T \neq 0$$

Otherwise: Representation **overcomplete**

F6

# Moment Parameters: Examples

$$P(\mathbf{x}; \theta) = Z^{-1} e^{\theta^T \mathbf{f}(\mathbf{x})}$$

$\theta$	Natural parameters
$\mathbf{f}(\mathbf{x})$	Statistics, representation
$\mu = E_{\theta}[\mathbf{f}(\mathbf{x})]$	Moment parameters

Representation **minimal**: For every  $\mathbf{z} \neq \mathbf{0}$ , there is  $\mathbf{x}$ :

$$\mathbf{z}^T (\mathbf{f}(\mathbf{x})^T \mathbf{1})^T \neq 0$$

Otherwise: Representation **overcomplete**

- ① Multinomial on graph with cliques  $C_j$

F6b

Convenient overcomplete representation: Components of  $\mathbf{f}(\mathbf{x})$ :

Indicators on cliques  $C_j$ , indicators on intersections of cliques, indicators on intersections of cliques, intersections, ...

Equality constraints for  $\mu$ :

- Consistency on nonempty intersections
- Sum to one on smallest intersections

# Moment Parameters: Examples

$$P(\mathbf{x}; \boldsymbol{\theta}) = Z^{-1} e^{\boldsymbol{\theta}^T \mathbf{f}(\mathbf{x})}$$

$\boldsymbol{\theta}$	Natural parameters
$\mathbf{f}(\mathbf{x})$	Statistics, representation
$\boldsymbol{\mu} = E_{\boldsymbol{\theta}}[\mathbf{f}(\mathbf{x})]$	Moment parameters

Representation **minimal**: For every  $\mathbf{z} \neq \mathbf{0}$ , there is  $\mathbf{x}$ :

$$\mathbf{z}^T (\mathbf{f}(\mathbf{x})^T \mathbf{1})^T \neq 0$$

Otherwise: Representation **overcomplete**

## 2 Gaussian MRF

Overcomplete representation:

F6c

$$\mathbf{f}(\mathbf{x}) = \begin{pmatrix} \mathbf{x} \\ \text{vec}(-\mathbf{x}\mathbf{x}^T/2) \end{pmatrix}, \quad \boldsymbol{\theta} = \begin{pmatrix} \mathbf{r} \\ \text{vec}(\mathbf{A}) \end{pmatrix}$$

Not minimal:  $\mathbf{A}$  symmetric.  $\{ij\} \notin E \rightarrow a_{ij} = a_{ji} = 0$ .

# Variational Formulation of Bayesian Inference

$$\log Z = \sup_Q \left\{ \boldsymbol{\theta}^T \mathbb{E}_Q[\mathbf{f}(\mathbf{x})] + H[Q(\mathbf{x})] \right\}, \quad \mathbf{f}(\mathbf{x}) = [\mathbf{f}_j(\mathbf{x}_{C_j})]$$

- Transform to moment parameters

# Variational Formulation of Bayesian Inference

$$\log Z = \sup_{\boldsymbol{\mu} \in \mathcal{M}} \left\{ \boldsymbol{\theta}^T \boldsymbol{\mu} + H[Q(\mathbf{x})] \right\}, \quad \mu_j = E_Q[\mathbf{f}_j(\mathbf{x}_{C_j})]$$

- Transform to moment parameters

$$\mathcal{M} = \left\{ (\boldsymbol{\mu}_j) \mid \mu_j = E_Q[\mathbf{f}_j(\mathbf{x}_{C_j})] \text{ for some } Q(\mathbf{x}) \right\}$$

⇒ **Marginal polytope**

F7

- What about the entropy?

# Variational Formulation of Bayesian Inference

$$\log Z = \sup_{\boldsymbol{\mu} \in \mathcal{M}} \left\{ \boldsymbol{\theta}^T \boldsymbol{\mu} + H[\boldsymbol{\mu}] \right\}$$

- Transform to moment parameters

$$\mathcal{M} = \left\{ (\boldsymbol{\mu}_j) \mid \boldsymbol{\mu}_j = \mathbb{E}_Q[\mathbf{f}_j(\mathbf{x}_{C_j})] \text{ for some } Q(\mathbf{x}) \right\}$$

⇒ **Marginal polytope**

- What about the entropy?  
 $\boldsymbol{\mu} \in \mathcal{M} \leftrightarrow Q(\mathbf{x})$  unique: Entropy **is** function of  $\boldsymbol{\mu}$
- Point of this exercise:  $\mathcal{M}$  convex set of vectors, more useful relaxation target than set of distributions

# Variational Formulation of Bayesian Inference

$$\log Z = \sup_{\boldsymbol{\mu} \in \mathcal{M}} \left\{ \boldsymbol{\theta}^T \boldsymbol{\mu} + H[\boldsymbol{\mu}] \right\}$$

- Transform to moment parameters

$$\mathcal{M} = \left\{ (\boldsymbol{\mu}_j) \mid \boldsymbol{\mu}_j = \mathbb{E}_Q[\mathbf{f}_j(\mathbf{x}_{C_j})] \text{ for some } Q(\mathbf{x}) \right\}$$

⇒ **Marginal polytope**

- What about the entropy?  
 $\boldsymbol{\mu} \in \mathcal{M} \leftrightarrow Q(\mathbf{x})$  unique: Entropy **is** function of  $\boldsymbol{\mu}$
- Point of this exercise:  $\mathcal{M}$  convex set of vectors, more useful relaxation target than set of distributions
- Close now: Exponential families, Fenchel duality, maximum entropy. Full story:

Wainwright, Jordan: Graphical Models, Exponential Families, and Variational Inference  
 Foundations and Trends in Machine Learning, 1(1–2), pp. 1–305

# Bayesian Inference is Convex Optimization

$$\log Z = \sup_{\mu \in \mathcal{M}} \left\{ \theta^T \mu + H[\mu] \right\}$$

- Marginal polytope  $\mathcal{M}$ : Convex set

F8



# Bayesian Inference is Convex Optimization

$$\log Z = \underbrace{\sup_{\mu \in \mathcal{M}}}_{\mathcal{M} \text{ convex}} \left\{ \theta^T \mu + H[\mu] \right\}$$

- Marginal polytope  $\mathcal{M}$ : Convex set
- Entropy  $\mu \mapsto H[\mu]$ : Concave function on  $\mathcal{M}$

F8b

# Bayesian Inference is Convex Optimization

$$\log Z = \underbrace{\sup_{\mu \in \mathcal{M}}}_{\mathcal{M} \text{ convex}} \underbrace{\left\{ \theta^T \mu + H[\mu] \right\}}_{\text{concave}}$$

- Marginal polytope  $\mathcal{M}$ : Convex set
- Entropy  $\mu \mapsto H[\mu]$ : Concave function on  $\mathcal{M}$
- Posterior: Unique solution to convex optimization problem

# Bayesian Inference is Convex Optimization

$$\log Z = \underbrace{\sup_{\mu \in \mathcal{M}}}_{\mathcal{M} \text{ convex}} \underbrace{\left\{ \theta^T \mu + H[\mu] \right\}}_{\text{concave}}$$

- Marginal polytope  $\mathcal{M}$ : Convex set
- Entropy  $\mu \mapsto H[\mu]$ : Concave function on  $\mathcal{M}$
- Posterior: Unique solution to convex optimization problem
- **Convex optimization can be intractable**

F8c

$\mathcal{M}$  can be hard to fence in  
 $\theta \leftrightarrow \mu$  can be hard to compute  
 $H[\mu]$  can be hard to compute

# Bayesian Inference is Convex Optimization

$$\log Z = \underbrace{\sup_{\mu \in \mathcal{M}}}_{\mathcal{M} \text{ convex}} \underbrace{\left\{ \theta^T \mu + H[\mu] \right\}}_{\text{concave}}$$

- Marginal polytope  $\mathcal{M}$ : Convex set
- Entropy  $\mu \mapsto H[\mu]$ : Concave function on  $\mathcal{M}$
- Posterior: Unique solution to convex optimization problem
- **Convex optimization can be intractable**
  - $\mathcal{M}$  can be hard to fence in
  - $\theta \leftrightarrow \mu$  can be hard to compute
  - $H[\mu]$  can be hard to compute
- Took some steps. But worth it:
  - Rich literature on **relaxations** of hard convex problems

# Variational Mean Field Revisited

$$\log Z = \sup_{\mu \in \mathcal{M}} \left\{ \theta^T \mu + H[\mu] \right\}$$

- Have to approximate  $\mathcal{M}$ ,  $H[\mu]$ . One way you already know . . .

# Variational Mean Field Revisited

$$\log Z = \sup_{\mu \in \mathcal{M}} \left\{ \theta^T \mu + H[\mu] \right\}$$

- Have to approximate  $\mathcal{M}$ ,  $H[\mu]$ . One way you already know ...

$$\mathcal{M} \supset \mathcal{M}_{\text{NMF}} := \left\{ \mu \mid \mu_{C_j}(\mathbf{x}_{C_j}) = \sum_{\mathbf{x}_{C_j}} \left( \prod_{i \in C_j} Q(x_i) \right) \mathbf{f}_j(\mathbf{x}_{C_j}) \right\}$$

**Inner approximation**, induced by factorized distributions

- Entropy decomposes just as distribution:  $H[\mu] = \sum_i H[\mu_i]$

$$\log Z \geq \sup_{\mu \in \mathcal{M}_{\text{NMF}}} \left\{ \theta^T \mu + \sum_i H[\mu_i] \right\}$$

# Variational Mean Field Revisited

$$\log Z = \sup_{\mu \in \mathcal{M}} \left\{ \theta^T \mu + H[\mu] \right\}$$

- Have to approximate  $\mathcal{M}$ ,  $H[\mu]$ . One way you already know ...

$$\mathcal{M} \supset \mathcal{M}_{\text{NMF}} := \left\{ \mu \mid \mu_{C_j}(\mathbf{x}_{C_j}) = \sum_{\mathbf{x}_{C_j}} \left( \prod_{i \in C_j} Q(x_i) \right) \mathbf{f}_j(\mathbf{x}_{C_j}) \right\}$$

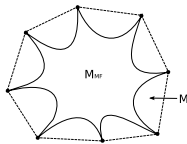
**Inner approximation**, induced by factorized distributions

- Entropy decomposes just as distribution:  $H[\mu] = \sum_i H[\mu_i]$

$$\log Z \geq \sup_{\mu \in \mathcal{M}_{\text{NMF}}} \left\{ \theta^T \mu + \sum_i H[\mu_i] \right\}$$

- Non-convex relaxation:  $\mathcal{M}_{\text{NMF}}$  **not convex**

F9



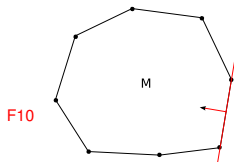
# The Marginal Polytope

$$\mathcal{M} = \left\{ (\boldsymbol{\mu}_j) \mid \mu_j = \mathbb{E}_Q[\mathbf{f}_j(\mathbf{x}_{C_j})] \text{ for some } Q(\mathbf{x}) \right\}$$

Multinomial on graph. Minimal representation.

- $\mathcal{M}$  convex polytope: Described by finite number inequalities.

Complexity of  $\mathcal{M}$ : Number of inequalities





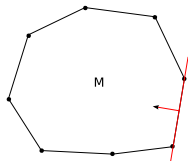
# The Marginal Polytope

$$\mathcal{M} = \left\{ (\boldsymbol{\mu}_j) \mid \mu_j = \mathbb{E}_Q[\mathbf{f}_j(\mathbf{x}_{C_j})] \text{ for some } Q(\mathbf{x}) \right\}$$

Multinomial on graph. Minimal representation.

- $\mathcal{M}$  convex polytope: Described by finite number inequalities.

Complexity of  $\mathcal{M}$ : Number of inequalities



- Complexity of  $\mathcal{M} \rightarrow$  complexity of exact inference [we'll see why]
- $\mathcal{G}$  tree:  $\mathcal{M}$  described by  $O(n)$  inequalities [next lecture]

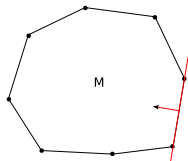
# The Marginal Polytope

$$\mathcal{M} = \left\{ (\boldsymbol{\mu}_j) \mid \mu_j = \mathbb{E}_Q[\mathbf{f}_j(\mathbf{x}_{C_j})] \text{ for some } Q(\mathbf{x}) \right\}$$

Multinomial on graph. Minimal representation.

- $\mathcal{M}$  convex polytope: Described by finite number inequalities.

Complexity of  $\mathcal{M}$ : Number of inequalities



- Complexity of  $\mathcal{M} \rightarrow$  complexity of exact inference [we'll see why]
- $\mathcal{G}$  tree:  $\mathcal{M}$  described by  $O(n)$  inequalities [next lecture]
- Many graphs  $\mathcal{G}$  with cycles:  $\mathcal{M}$  polytope description provably hard (poly( $n$ ) inequalities would imply P=NP)

# The Marginal Polytope

$$\mathcal{M} = \left\{ (\boldsymbol{\mu}_j) \mid \boldsymbol{\mu}_j = \mathbb{E}_Q[\mathbf{f}_j(\mathbf{x}_{C_j})] \text{ for some } Q(\mathbf{x}) \right\}$$

Gaussian MRF. Minimal representation (upper triangle of  $\mathbf{A}$ ).

- $\mathcal{M}$  exactly characterized by  $\boldsymbol{\Sigma} = \mathbf{A}^{-1} \succ \mathbf{0}$ .  
Convex cone (not polytope): Tractable to describe
- $\mathcal{G}$  tree:  $\mathcal{M}$  described by  $O(n)$  inequalities
- General sparse  $\mathcal{G}$ : Approximate inference still of interest, if exact cost  $O(n^3)$  too high

F10b

# Wrap-Up

- Structured Mean Field:  $Q(\mathbf{x})$  product of tractable, disjoint factors
- Variational Bayes: Another name for structured mean field
- Bayesian (marginal) inference is a convex optimization problem
- Variational approximations: Inner / outer bounds to marginal polytope