

ECE 6504: Advanced Topics in Machine Learning

Probabilistic Graphical Models and Large-Scale Learning

Topics

- Markov Random Fields: **MAP** Inference
 - Max-Product Message Passing
 - Integer Programming, LP formulation
 - Dual Decomposition

Readings: KF 13.1-5, Barber 5.1,28.9

Dhruv Batra

Virginia Tech

Administrativa

- HW1 Solutions
 - Released
 - Grades almost done too
- Project Presentations
 - When: April 22, 24
 - Where: in class
 - 5 min talk
 - Main results
 - Semester completion 2 weeks out from that point so nearly finished results expected
 - Slides due: April 21 11:55pm



Recap of Last Time

Message Passing

- Variables/Factors “talk” to each other via messages:

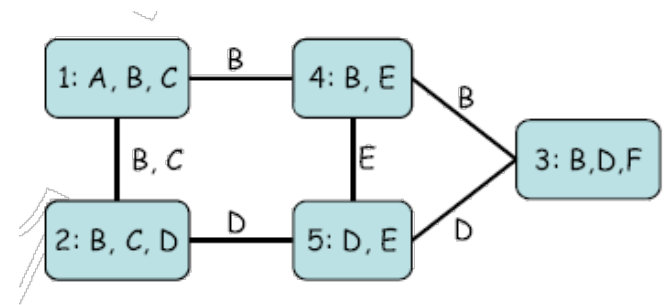
“I (variable X_3) think that you (variable X_2):
belong to state 1 with confidence 0.4
belong to state 2 with confidence 10
belong to state 3 with confidence 1.5”



Generalized BP

- Initialization:

- Assign each factor ϕ to a cluster $\alpha(\phi)$, $\text{Scope}[\phi] \subseteq \mathbf{C}_{\alpha(\phi)}$
- Initialize cluster: $\psi_i^0(\mathbf{C}_i) \propto \prod_{\phi: \alpha(\phi)=i} \phi$
- Initialize messages: $\delta_{j \rightarrow i} = 1$



- While not converged, send messages:

$$\delta_{i \rightarrow j}(\mathbf{S}_{ij}) \propto \sum_{\mathbf{C}_i - \mathbf{S}_{ij}} \psi_i^0(\mathbf{C}_i) \prod_{k \in \mathcal{N}(i) - j} \delta_{k \rightarrow i}(\mathbf{S}_{ik})$$

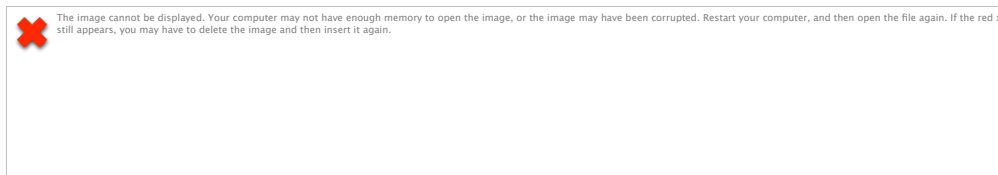
- Belief:

- On board

What is Variational Inference?

- A class of methods for approximate inference
 - And parameter learning
 - And approximating integrals basically..
- Key idea
 - Reality is complex
 - Instead of performing approximate computation in something complex
 - Can we perform exact computation in something “simple”?
 - Just need to make sure the simple thing is “close” to the complex thing.
- Key Problems
 - What is close?
 - How do we measure closeness when we can't perform operations on the complex thing?

Variational Approximate Inference



- Choose a family of approximating distributions which is tractable. The simplest [Mean Field] Approximation:

$$q(x) = \prod_{s \in \mathcal{V}} q_s(x_s)$$

- Measure the quality of approximations. Two possibilities:

$$D(p \parallel q) = \sum_x p(x) \log \frac{p(x)}{q(x)} \quad D(q \parallel p) = \sum_x q(x) \log \frac{q(x)}{p(x)}$$

- Find the approximation minimizing this distance

D(p||q) for mean field – KL the right way

- D(p||q)=

- Trivially minimized by setting $q_i(x_i) = p_i(x_i)$
- Doesn't provide a computational method...

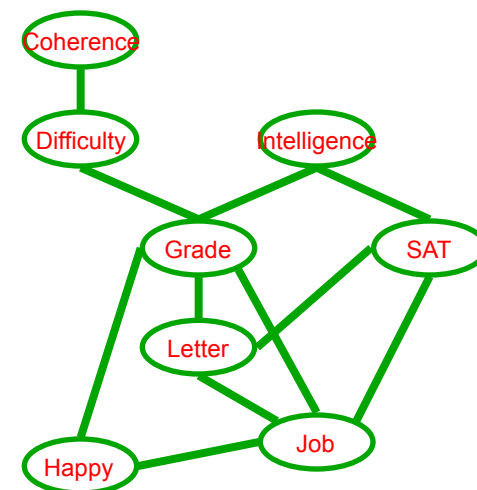
D(q||p) for mean field – KL the reverse direction

- D(q||p)=

$$D(q||p) = \sum_x q(x) \log q(x) - \sum_x q(x) \log p(x)$$

Reverse KL & The Partition Function

- $D(q||p)$:
 - p is Markov net P_F



- **Theorem:** $\log Z = F[p, q] + D(q||p)$

- Where “Gibbs Free Energy”:

$$F[p, q] = H_q(\mathcal{X}) + \mathbb{E}_q \left[\sum_c \log \psi_c(X_c) \right]$$

$$= H_q(\mathcal{X}) + \mathbb{E}_q [\text{Score}(\mathcal{X})]$$

$$= H_q(\mathcal{X}) + \sum_c \sum_{x_c} q(x_c) \theta(x_c)$$

Understanding Reverse KL, Free Energy & The Partition Function

$$\log Z = F[p, q] + D(q||p) \qquad F[p, q] = H_q(\mathcal{X}) + \mathbb{E}_q \left[\sum_c \log \psi_c(X_c) \right]$$

- Maximizing Energy Functional \Leftrightarrow Minimizing Reverse KL
- **Theorem:** Energy Function is lower bound on partition function
 - Maximizing energy functional corresponds to search for tight lower bound on partition function

Mean Field Equations

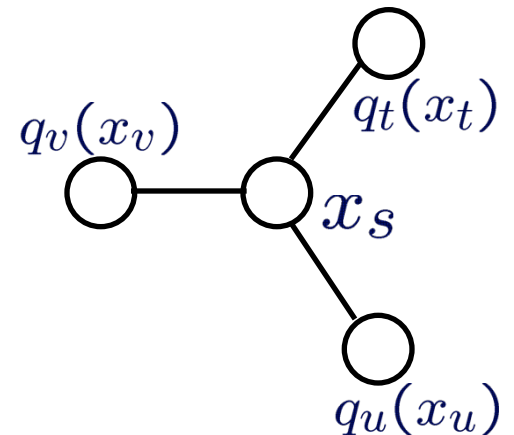
$$F[p, q] = H_q(\mathcal{X}) + \mathbb{E}_q \left[\sum_c \log \psi_c(X_c) \right]$$

$$H(q) = \sum_{s \in \mathcal{V}} H_s(q_s) = - \sum_{s \in \mathcal{V}} \sum_{x_s} q_s(x_s) \log q_s(x_s)$$

$$\sum_c \sum_{x_c} q_c(x_c) \theta(x_c) = \sum_i \sum_{x_i} q_i(x_i) \theta_i(x_i) + \sum_{(i,j) \in E} \sum_{x_i} \sum_{x_j} q_i(x_i) q_j(x_j) \theta_{ij}(x_i, x_j)$$

- Add Lagrange multipliers to enforce $\sum_{x_s} q_s(x_s) = 1$
- Taking derivatives and simplifying, we find a set of fixed point equations:

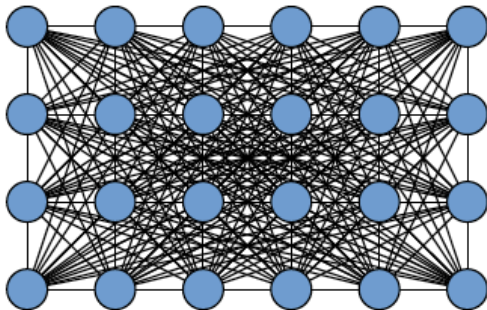
$$q_i(x_i) \propto \psi_i(x_i) \prod_{j \in N(i)} \exp \left\{ \sum_{x_j} \theta_{ij}(x_i, x_j) q_j(x_j) \right\}$$



- Updating one marginal at a time gives convergent coordinate descent

Fully connected CRF

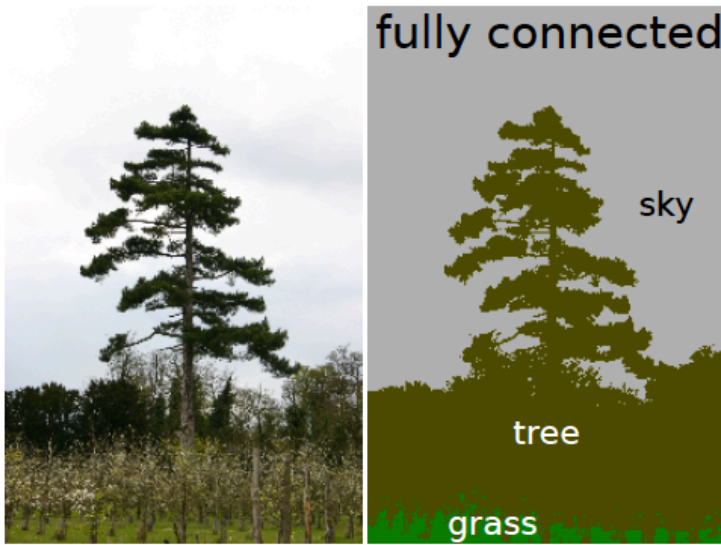
$$E(\mathbf{x}) = \sum_i \underbrace{\psi_u(x_i)}_{\text{unary term}} + \sum_i \sum_{j>i} \underbrace{\psi_p(x_i, x_j)}_{\text{pairwise term}}$$



- Every node is connected to every other node
 - ▶ Connections weighted differently

Fully connected CRF

$$E(\mathbf{x}) = \sum_i \underbrace{\psi_u(x_i)}_{\text{unary term}} + \sum_i \sum_{j>i} \underbrace{\psi_p(x_i, x_j)}_{\text{pairwise term}}$$



- Long-range interactions
- No more shrinking bias

What you need to know about variational methods

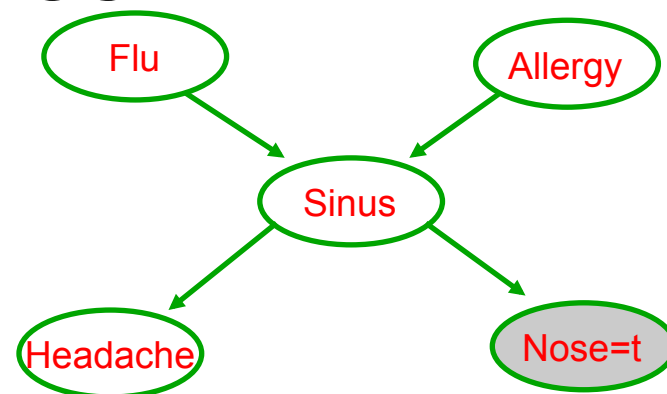
- Structured Variational method:
 - select a form for approximate distribution
 - minimize reverse KL
- Equivalent to maximizing energy functional
 - searching for a tight lower bound on the partition function
- Many possible models for Q :
 - independent (mean field)
 - structured as a Markov net
 - cluster variational
- Several subtleties outlined in the book

Plan for today

- MRF Inference
 - (Specialized) MAP Inference
 - Integer Programming Formulation
 - Linear Programming Relaxation
 - Dual Decomposition

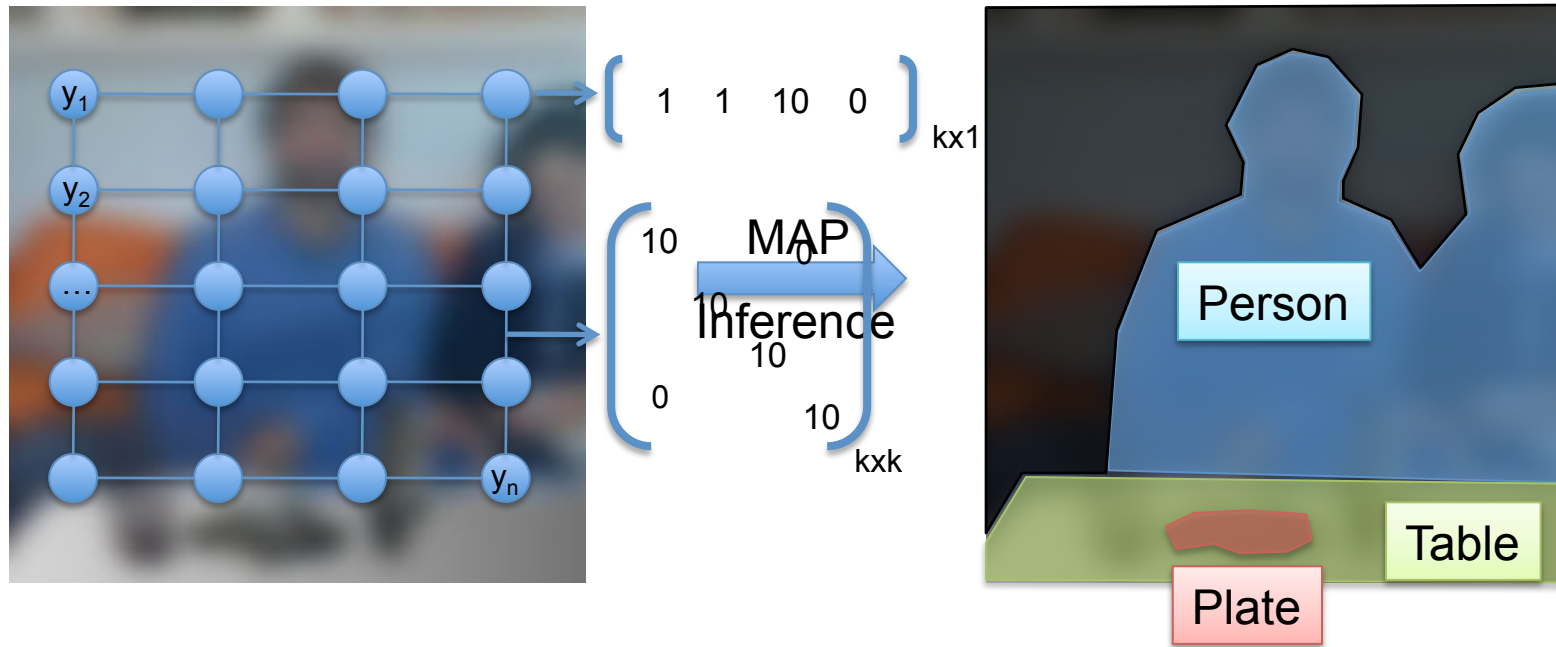
Possible Queries

- Evidence: $\mathbf{E}=\mathbf{e}$ (e.g. $N=t$)
- Query variables of interest \mathbf{Y}



- Conditional Probability: $P(\mathbf{Y} \mid \mathbf{E}=\mathbf{e})$
 - E.g. $P(F,A \mid N=t)$
 - **Special case:** Marginals $P(F)$
- Maximum a Posteriori: $\operatorname{argmax} P(\text{All variables} \mid \mathbf{E}=\mathbf{e})$
 - $\operatorname{argmax}_{\{f,a,s,h\}} P(f,a,s,h \mid N = t)$
- Marginal-MAP: $\operatorname{argmax}_y P(\mathbf{Y} \mid \mathbf{E}=\mathbf{e})$
 - $= \operatorname{argmax}_{\{y\}} \sum_o P(\mathbf{Y}=\mathbf{y}, \mathbf{O}=\mathbf{o} \mid \mathbf{E}=\mathbf{e})$

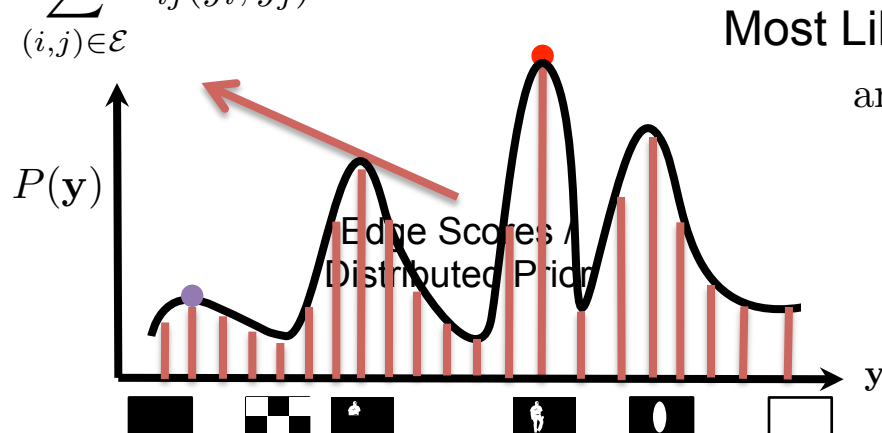
MAP Inference



$$S(\mathbf{y}) = \sum_{i \in \mathcal{V}} \theta_i(y_i) + \sum_{(i,j) \in \mathcal{E}} \theta_{ij}(y_i, y_j)$$

$$P(\mathbf{y}) = \frac{1}{Z} e^{S(\mathbf{y})}$$

Node Scores /
Local Rewards

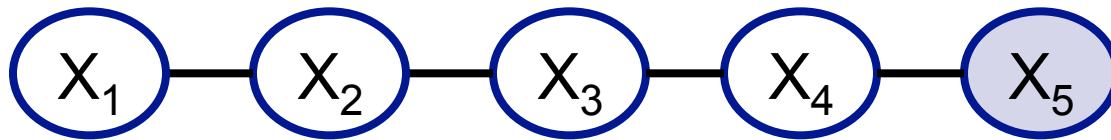


Most Likely Assignment

$$\operatorname{argmax}_{\mathbf{y}} S(\mathbf{y})$$

Example

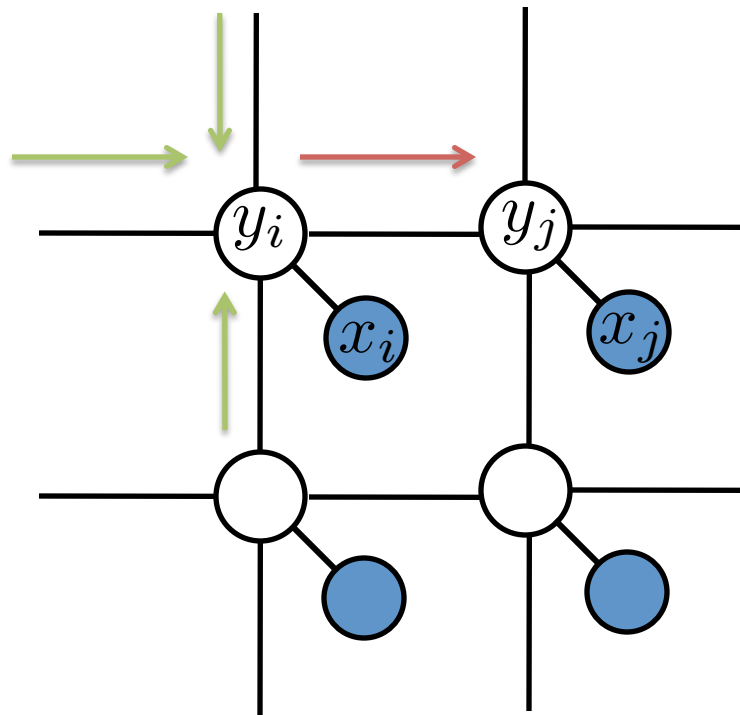
- Chain MRF



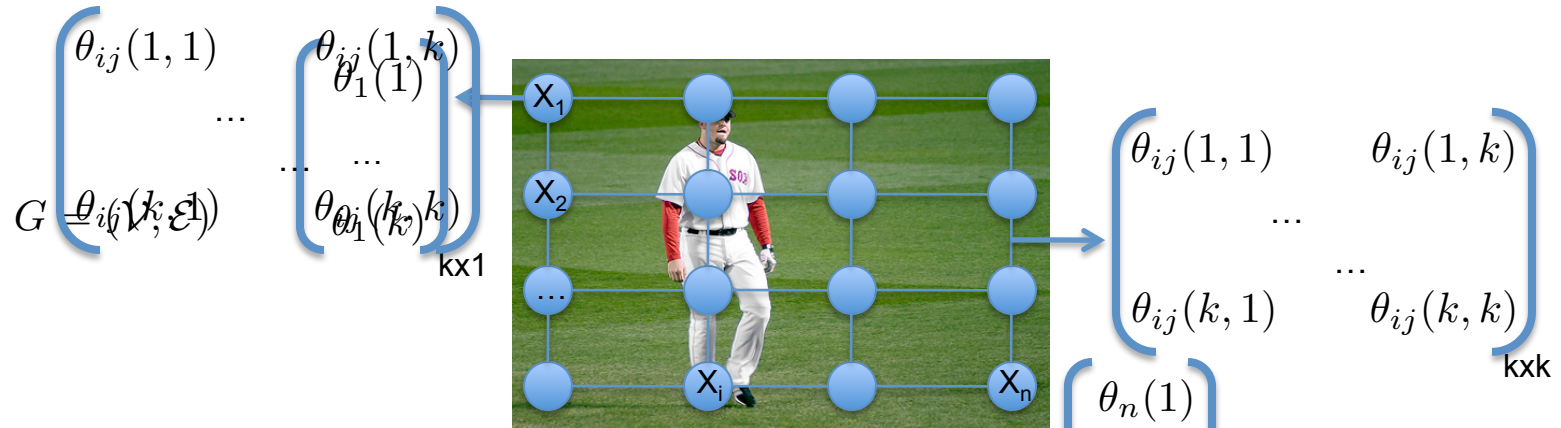
- Max-Product VE steps on board

Loopy BP on Pairwise Markov Nets

$$\overrightarrow{\delta}_{i \rightarrow j}(y_j) = \sum_{y_i} \phi_i(y_i) \phi_{ij}(y_i, y_j) \prod_{k \in \mathcal{N}(i) - j} \overrightarrow{\delta}_{k \rightarrow i}(y_i)$$

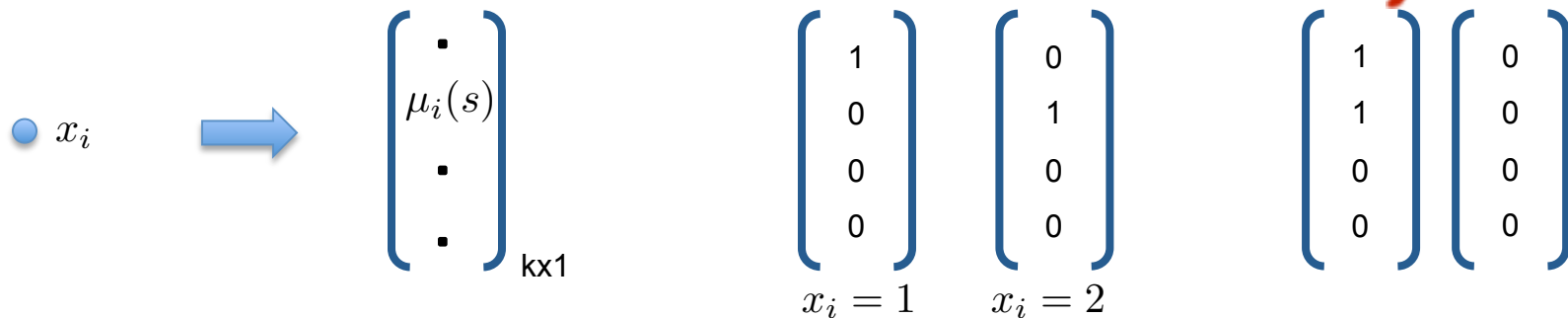


MAP in Pairwise MRFs



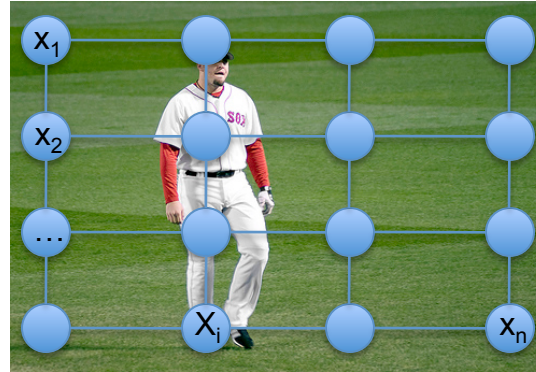
- Over-Complete Representation

$$\theta = \begin{matrix} \boxed{x_1} & \dots & \boxed{x_n} & \boxed{(x_1, x_2)} & \dots & \boxed{(x_{n-1}, x_n)} \\ \left[\begin{matrix} \theta_1(1) \dots \theta_1(k) & \theta_n(1) \dots \theta_n(k) & \theta_{12}(1,1) \dots \theta_{12}(k,k) & \theta_{n-1,n}(1,1) \dots \theta_{n-1,n}(k,k) \end{matrix} \right] \\ \mu_1(1) \dots \mu_1(k) & & \mu_n(1) \dots \mu_n(k) & & & \end{matrix}$$



MAP in Pairwise MRFs

$$G = (\mathcal{V}, \mathcal{E})$$



- Over-Complete Representation

$$\begin{array}{cccc}
 \boxed{x_1} & \cdots & \boxed{x_n} & \boxed{(x_1, x_2)} & \cdots & \boxed{(x_{n-1}, x_n)} \\
 \theta = & \left[\begin{array}{ccc} \theta_1(1) \dots \theta_1(k) & \theta_n(1) \dots \theta_n(k) & \theta_{12}(1,1) \dots \theta_{12}(k,k) \end{array} \right] & & \left[\begin{array}{c} \theta_{n-1,n}(1,1) \dots \theta_{n-1,n}(k,k) \end{array} \right] \\
 \mu_{\mathbf{x}} = & \left[\begin{array}{ccc} \mu_1(1) \dots \mu_1(k) & \mu_n(1) \dots \mu_n(k) & \mu_{12}(1,1) \dots \mu_{12}(k,k) \end{array} \right] & & \left[\begin{array}{c} \mu_{n-1,n}(1,1) \dots \mu_{n-1,n}(k,k) \end{array} \right] \\
 \begin{array}{c} \bullet x_i \\ | \\ \bullet x_j \end{array} & \xrightarrow{\quad} & \begin{array}{c} \left[\begin{array}{c} \vdots \\ \mu_{ij}(s,t) \end{array} \right]_{k^2 \times 1} \\ S(\mathbf{x}) = \theta \cdot \mu_{\mathbf{x}} \end{array} & & \begin{array}{c} \left[\begin{array}{c} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \right] \begin{array}{l} x_i = 1 \\ x_j = 1 \end{array} \\ \left[\begin{array}{c} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \right] \begin{array}{l} x_i = 1 \\ x_j = 2 \end{array} \end{array}
 \end{array}$$

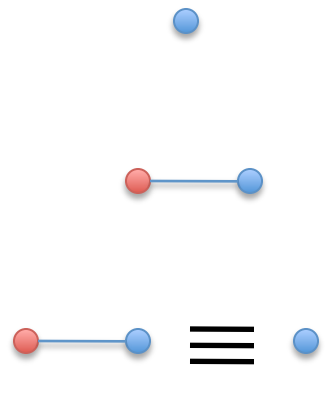
MAP in Pairwise MRFs

- Integer Program

$$\max_{\mu} \theta^T \mu$$

$$\left. \begin{aligned} \mu_i(s) &\in \{0, 1\} \\ \mu_{ij}(s, t) &\in \{0, 1\} \end{aligned} \right\} \leftarrow \text{Indicator Variables}$$

$$\left. \begin{aligned} \sum_s \mu_i(s) &= 1 \\ \sum_{s,t} \mu_{i,j}(s, t) &= 1 \end{aligned} \right\} \leftarrow \text{Unique Label}$$

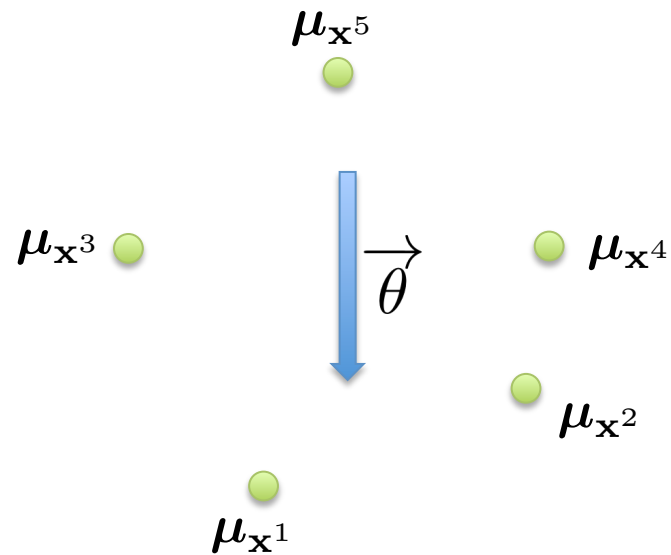


$$\left. \begin{aligned} \sum_s \mu_{ij}(s, t) &= \mu_j(t) \end{aligned} \right\} \leftarrow \text{Consistent Assignments}$$

MAP in Pairwise MRFs

- MAP Integer Program

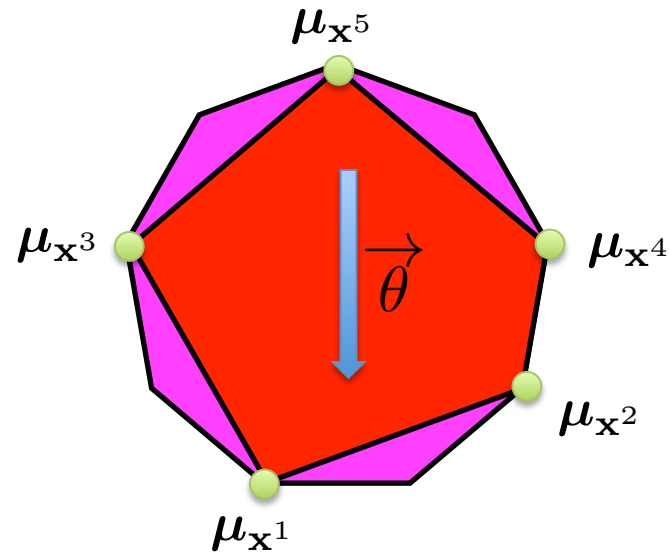
$$\begin{aligned} \max_{\mu} \quad & \theta^T \mu \\ \text{s.t.} \quad & A\mu = b \\ & \mu(\cdot) \in \{0, 1\} \end{aligned}$$



MAP in Pairwise MRFs

- MAP Linear Program

$$\begin{aligned} \max_{\mu} \quad & \theta^T \mu \\ \text{s.t.} \quad & A\mu = b \\ & \mu(\cdot) \in [0, 1] \end{aligned}$$



- Properties
 - If LP-opt is integral, MAP is found
 - LP always integral for trees
 - Efficient message-passing schemes for solving LP

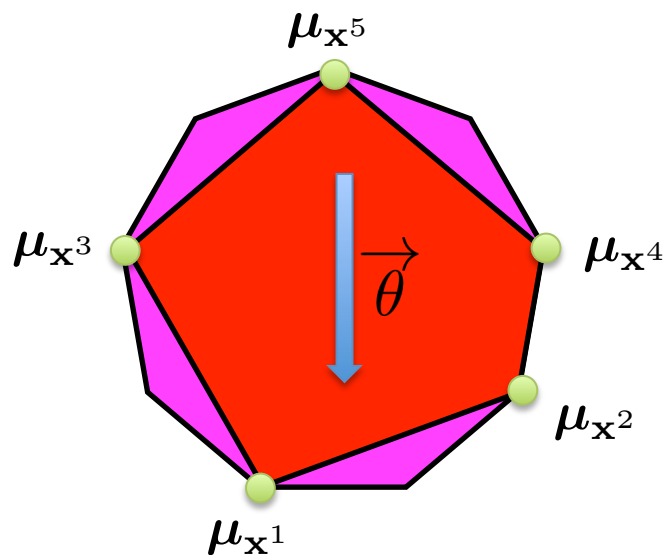
MAP in Pairwise MRFs

- Compare MAP LP to Variational Inference
 - On board
 - Difference in entropy term (objective)
 - Family of Q distributions

MAP in Pairwise MRFs

- MAP Linear Program

$$\begin{aligned} \max_{\mu} \quad & \theta^T \mu \\ \text{s.t.} \quad & A\mu = b \\ & \mu(\cdot) \in [0, 1] \end{aligned}$$



$$A = \begin{pmatrix} \\ \\ \\ \\ \end{pmatrix} \begin{matrix} \text{ } \\ \text{ } \\ \text{ } \\ \text{ } \\ \text{ } \end{matrix} \begin{matrix} \\ \\ \\ \\ \end{matrix}$$

$O(|\mathcal{E}|)$ (vertical dimension)

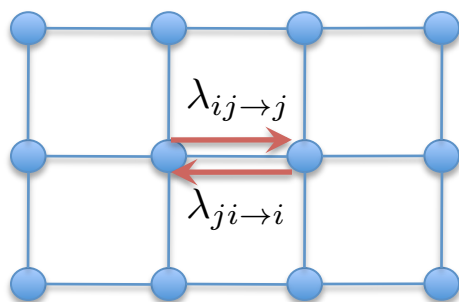
$O(|\mathcal{E}|)$ (horizontal dimension)

Off-the-shelf solvers
CPLEX
Mosek
etc

✗

LP Relaxation

- Block Co-ordinate / Sub-gradient Descent on Dual

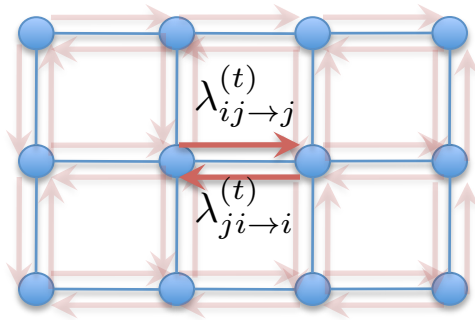


$$A = \left[\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right] \begin{array}{l} \updownarrow \\ O(|\mathcal{E}|) \end{array}$$

$\leftarrow \text{---} \rightarrow$
 $O(|\mathcal{E}|)$

LP Relaxation

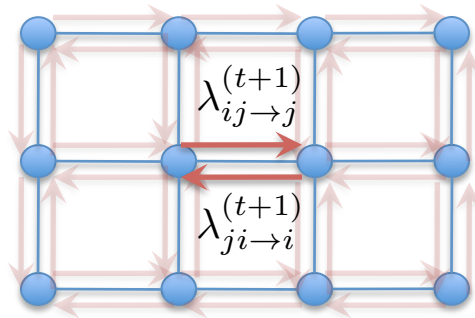
- Block Co-ordinate / Sub-gradient Descent on Dual



$$A = \left(\begin{array}{c} \\ \\ \\ \end{array} \right) \begin{array}{l} \updownarrow \\ O(|\mathcal{E}|) \end{array}$$
$$\begin{array}{c} \leftarrow \rightleftarrows \\ O(|\mathcal{E}|) \end{array}$$

LP Relaxation

- Block Co-ordinate / Sub-gradient Descent on Dual



Distributed Message-Passing

Still inefficient!

$$A = \left[\begin{array}{c} \\ \\ \end{array} \right] \begin{array}{c} \updownarrow \\ O(|\mathcal{E}|) \end{array}$$
$$\begin{array}{c} \leftarrow \\ O(|\mathcal{E}|) \rightarrow \end{array}$$