

Graphical Models

Lecture 13:

Inference as Optimization, Mean Field

Andrew McCallum
mccallum@cs.umass.edu

Thanks to Noah Smith and Carlos Guestrin for some slide materials.

Administration

- Homework 3 source code due Thursday

Toward Approximate Inference

- Exact inference is not always tractable.
 - NP hard in general!
- Efficient, principled approximations are incredibly useful in practice.
 - There are a lot of them!
 - We'll cover
 - mean field variational inference
 - loopy belief propagation (also variational)
 - Markov chain Monte Carlo

General Approach

- P is a “hard” distribution.
- Pick a class of “easy” distributions \mathcal{Q} over the same random variables.
- Different views:
 - Optimization: find $Q (\in \mathcal{Q})$ to minimize the “distance” between P and Q.
 - Fixed-points
 - Some form of message passing or updating

KL Divergence (Relative Entropy)

$$D(P\|Q) = \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} P(\mathbf{X} = \mathbf{x}) \log \frac{P(\mathbf{X} = \mathbf{x})}{Q(\mathbf{X} = \mathbf{x})}$$

- A measurement of “distance” between two distributions.
 - Not symmetric.
- For exponential family P and any Q:

$$D(Q\|P_{\boldsymbol{\theta}}) = -H(Q(\mathbf{X})) - \mathbf{t}(\boldsymbol{\theta})^{\top} \mathbb{E}_Q[\boldsymbol{\tau}(\mathbf{X})] + \ln Z(\boldsymbol{\theta})$$

natural parameter

sufficient statistics

Lecture 10

Projections

- Given a distribution P and an exponential family \mathcal{Q} , find the distribution from \mathcal{Q} that is closest to P .

– I-projection (information projection):

$$\arg \min_{Q \in \mathcal{Q}} D(Q \| P)$$

– M-projection (moment projection):

$$\arg \min_{Q \in \mathcal{Q}} D(P \| Q)$$

Projections

- Given a distribution P and an exponential family \mathcal{Q} , find the distribution from \mathcal{Q} that is closest to P .

- I-projection (information projection):

$$\arg \min_{Q \in \mathcal{Q}} D(Q \| P)$$

- M-projection (moment projection):

$$\arg \min_{Q \in \mathcal{Q}} D(P \| Q)$$

minimize number of bits lost when coding a true distribution P using approximate Q ; but requires inference in P (expectations under P).

Projections

- Given a distribution P and an exponential family \mathcal{Q} , find the distribution from \mathcal{Q} that is closest to P .

- I-projection (information projection):

$$\arg \min_{Q \in \mathcal{Q}} D(Q \| P)$$

easier
(as we will see)

- M-projection (moment projection):

$$\arg \min_{Q \in \mathcal{Q}} D(P \| Q)$$

Claim: We (Almost) Already Do I-Projection!

- Let the class \mathcal{Q} be defined by $\{\beta_i\}$ and $\{\mu_{i,j}\}$.

$$\arg \min_{Q \in \mathcal{Q}} D(Q \| P)$$

subject to calibration constraints:

$$\forall (i, j) \in \mathcal{E}, \forall \mathbf{s}_{i,j} \in \text{Val}(\mathbf{S}_{i,j}), \quad \mu_{i,j}(\mathbf{s}_{i,j}) = \sum_{\mathbf{c}_i \setminus \mathbf{S}_{i,j}} \beta_i(\mathbf{c}_i)$$
$$\forall i \in \mathcal{V}, \quad \sum_{\mathbf{c}_i} \beta_i(\mathbf{c}_i) = 1$$

Claim: We (Almost) Already Do I-Projection!

- Let the class \mathcal{Q} be defined by $\{\beta_i\}$ and $\{\mu_{i,j}\}$.

$$\arg \min_{Q \in \mathcal{Q}} D(Q \| P)$$

stays zero throughout the BU message passing algorithm

subject to calibration constraints:

$$\forall (i, j) \in \mathcal{E}, \forall \mathbf{s}_{i,j} \in \text{Val}(\mathbf{S}_{i,j}), \quad \mu_{i,j}(\mathbf{s}_{i,j}) = \sum_{\mathbf{c}_i \setminus \mathbf{S}_{i,j}} \beta_i(\mathbf{c}_i)$$
$$\forall i \in \mathcal{V}, \quad \sum_{\mathbf{c}_i} \beta_i(\mathbf{c}_i) = 1$$

Claim: We (Almost) Already Do I-Projection!

- Let the class \mathcal{Q} be defined by $\{\beta_i\}$ and $\{\mu_{i,j}\}$.

$$\arg \min_{Q \in \mathcal{Q}} D(Q \| P)$$

calibration;
achieved at
convergence

subject to calibration constraints:

$$\forall (i, j) \in \mathcal{E}, \forall \mathbf{s}_{i,j} \in \text{Val}(\mathbf{S}_{i,j}),$$

$$\mu_{i,j}(\mathbf{s}_{i,j}) = \sum_{\mathbf{c}_i \setminus \mathbf{S}_{i,j}} \beta_i(\mathbf{c}_i)$$

$$\forall i \in \mathcal{V}, \sum_{\mathbf{c}_i} \beta_i(\mathbf{c}_i) = 1$$

Claim

- If the clique tree structure is an I-map for P , then there is a unique solution to this problem, found using the message passing algorithms we've already seen.

Alternative derivation for clique-tree belief update rules!

$$\arg \min_{Q \in \mathcal{Q}} D(Q \| P)$$

such that

$$\forall (i, j) \in \mathcal{E}, \forall \mathbf{s}_{i,j} \in \text{Val}(\mathbf{S}_{i,j}), \quad \mu_{i,j}(\mathbf{s}_{i,j}) = \sum_{\mathbf{c}_i \setminus \mathbf{S}_{i,j}} \beta_i(\mathbf{c}_i)$$

$$\forall i \in \mathcal{V}, \quad \sum_{\mathbf{c}_i} \beta_i(\mathbf{c}_i) = 1$$

I-projection & Helmholtz Free Energy

- (Board work)
- Derive using just definitions of KL divergence and Gibbs distribution.
- Energy functional as lower bound on log partition function.

More General Goal

- Define “easy” family of distributions \mathcal{Q} .
- Minimize $D(Q || P)$
- Assume a factorized form for Q that offers convenient structure.

$$\begin{aligned} P(\mathbf{X}) &= \frac{U(\mathbf{X})}{Z} \\ &= \frac{1}{Z} \prod_{\phi \in \Phi} \phi(\text{Scope}(\mathbf{X}; \phi)) \end{aligned}$$

$$\begin{aligned} D(Q||P) &= \mathbb{E}_Q[\log Q] - \mathbb{E}_Q[\log P] \\ &= -H_Q - (\mathbb{E}_Q[\log U] - \log Z) \\ &= -H_Q - \left(\sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi] - \log Z \right) \end{aligned}$$

$$\log Z = D(Q||P) + H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi]$$

More General Goal

- Define “easy” family of distributions \mathcal{Q} .
- Minimize $D(Q || P)$
- Assume a factorized form for Q that offers convenient structure.

$$\begin{aligned}
 D(Q||P) &= \mathbb{E}_Q[\log Q] - \mathbb{E}_Q[\log P] \\
 &= -H_Q - (\mathbb{E}_Q[\log U] - \log Z) \\
 &= -H_Q - \left(\sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi] - \log Z \right)
 \end{aligned}$$

constant in Q

$$\boxed{\log Z} = \underbrace{\boxed{D(Q||P)}}_{\geq 0} + \underbrace{\boxed{H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi]}}_{\text{energy functional}}$$

More General Goal

- Maximize a lower bound on log partition function!
 - In a directed model with evidence, Z is the posterior we care about, $P(\text{Evidence})$.

$$\begin{aligned}D(Q\|P) &= \mathbb{E}_Q[\log Q] - \mathbb{E}_Q[\log P] \\ &= -H_Q - (\mathbb{E}_Q[\log U] - \log Z) \\ &= -H_Q - \left(\sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi] - \log Z \right) \\ \log Z &= D(Q\|P) + H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi]\end{aligned}$$

Variational Methods

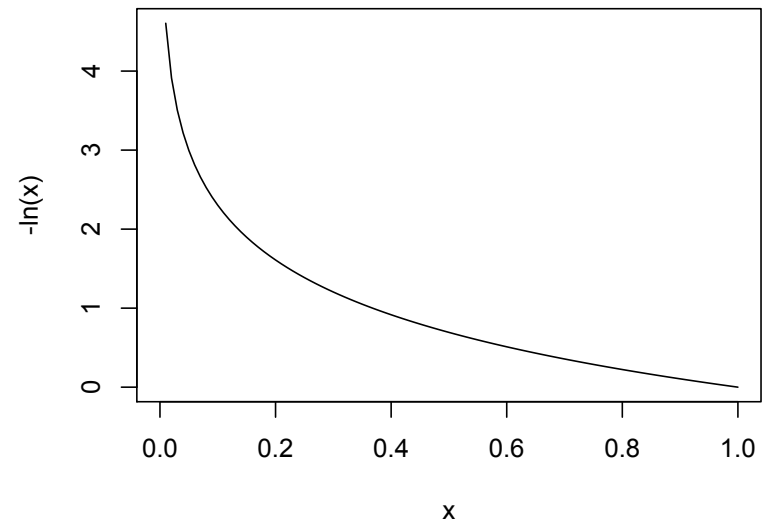
- The free variables here will be the parameters of Q , which comes from \mathcal{Q} .
- **Variational method:** optimize the energy functional.
 - Every element of \mathcal{Q} is an approximate solution.
 - We try to find the best one.

$$\max_{Q \in \mathcal{Q}} H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi] \equiv \min_{Q \in \mathcal{Q}} D(Q \| P)$$

Tangent: Variational Methods

- This is a simple example.
- For any λ and any x :

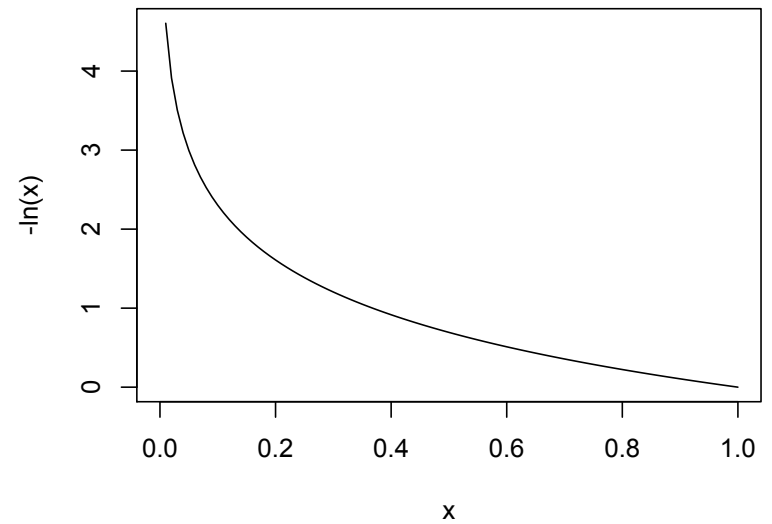
$$-\ln(x) \geq \underbrace{-\lambda x + \ln(\lambda) + 1}_{\text{family of functions } g_\lambda(x)}$$



Tangent: Variational Methods

- This is a simple example.
- For any λ and any x :

$$-\ln(x) \geq -\lambda x + \ln(\lambda) + 1$$

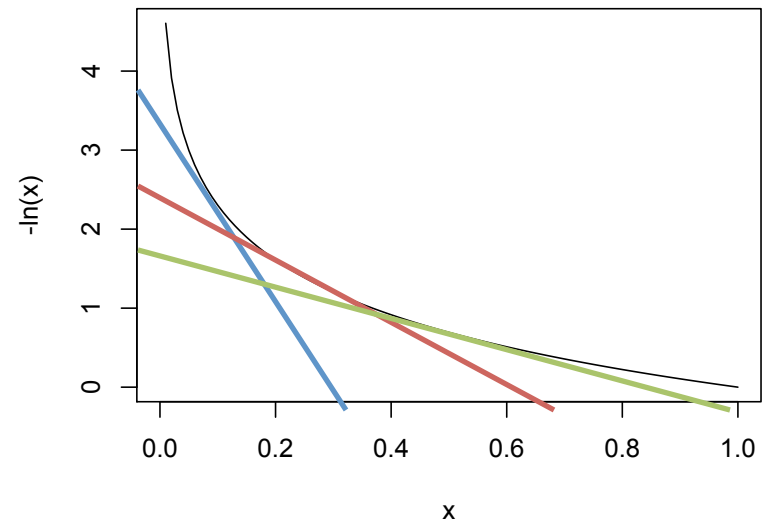


- Further, for any x , there is some λ where the bound is tight.
 - λ is called a **variational parameter**.

Tangent: Variational Methods

- This is a simple example.
- For any λ and any x :

$$-\ln(x) \geq -\lambda x + \ln(\lambda) + 1$$

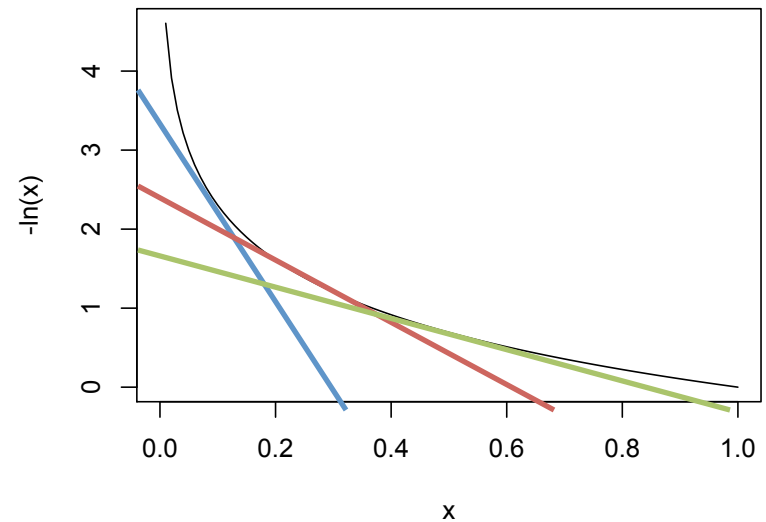


- Further, for any x , there is some λ where the bound is tight.
 - λ is called a **variational parameter**.

Tangent: Variational Methods

- This is a simple example.
- For any λ and any x :

$$-\ln(x) \geq -\lambda x + \ln(\lambda) + 1$$



- Further, for any x , there is some λ where the bound is tight.
 - λ is called a **variational parameter**.
- For us, $\log(Z)$ is like $-\ln(x)$, and Q is like λ .

the complex thing we
are trying to approximate

the parameters of the
simpler approximation
we want to be close

Inference as Optimization: 3 flavors

- All optimizing a distance between Q and P ,
e.g. $\arg \min_{Q \in \mathcal{Q}} D(Q \| P)$
- **mean field**
 - exact energy functional, but restricted class of \mathcal{Q}
(simple factorization)
- **loopy belief propagation**
 - approximate energy functional
- **expectation propagation**
 - exact energy functional, but approximate messages
(relaxed consistency constraints on Q)

Structured Variational Approach

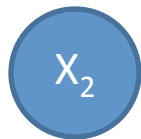
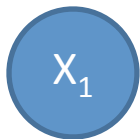
“Mean Field” “Structured Mean Field”

- Maximize the energy functional over a family \mathcal{Q} that is well-defined.
 - A graphical model!
 - Probably not an I-map for P . (Bound isn't tight.)
 - Simpler structures lead to easier inference.

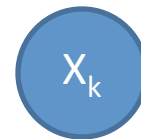
Mean Field

- The simplest of all possible families.

$$Q(\mathbf{X}) = \prod_i Q_i(X_i)$$



...



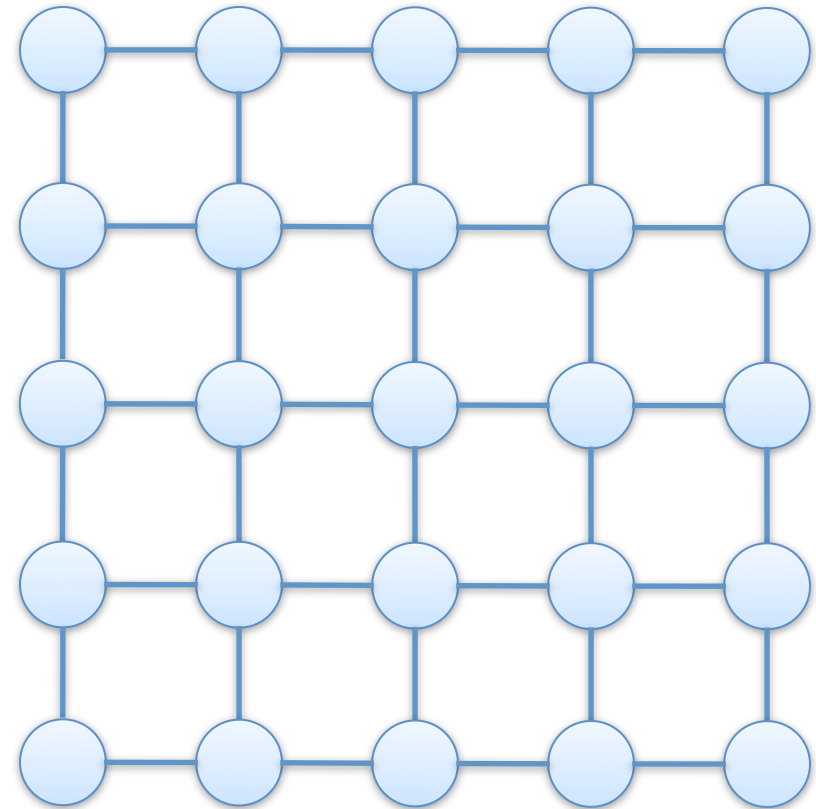
Pairwise Markov Network Example

- Classify each pixel as foreground or background.

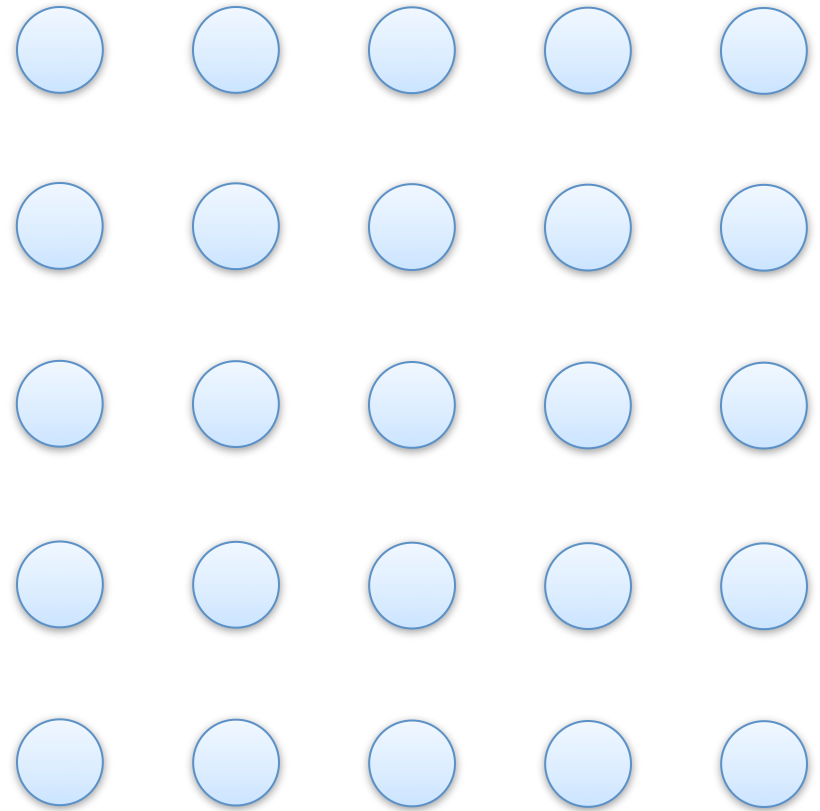
$$\phi_i(\text{fg}) = \exp \frac{-\|c_i - \mu_{\text{fg}}\|^2}{\sigma^2}$$

$$\phi_i(\text{bg}) = \exp \frac{-\|c_i - \mu_{\text{bg}}\|^2}{\sigma^2}$$

$$\phi_{i,j}(X_i, X_j) = \begin{cases} 10 & \text{if } X_i = X_j \\ 1 & \text{otherwise} \end{cases}$$



Example: Mean Field Q



Mean Field and the Energy Functional

$$\begin{aligned}
 & \max_{Q \in \mathcal{Q}} H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi] \\
 &= \max_{Q \in \mathcal{Q}} \sum_i H_{Q_i} + \sum_{\phi \in \Phi} \sum_{\mathbf{u} \in \text{Val}(\text{Scope}(\mathbf{X}; \phi))} Q(\mathbf{u}) \log \phi(\mathbf{u}) \\
 &= \max_{Q \in \mathcal{Q}} \sum_i H_{Q_i} + \sum_{\phi \in \Phi} \sum_{\mathbf{u} \in \text{Val}(\text{Scope}(\mathbf{X}; \phi))} \left(\prod_{X_i \in \text{Scope}(\mathbf{X}; \phi)} Q_i(u_i) \right) \log \phi(\mathbf{u}) \\
 &= \max_{Q \in \mathcal{Q}} - \sum_i \sum_{x \in \text{Val}(X_i)} Q_i(x_i) \log Q_i(x_i) \\
 &\quad + \sum_{\phi \in \Phi} \sum_{\mathbf{u} \in \text{Val}(\text{Scope}(\mathbf{X}; \phi))} \left(\prod_{X_i \in \text{Scope}(\mathbf{X}; \phi)} Q_i(u_i) \right) \log \phi(\mathbf{u})
 \end{aligned}$$

Optimization

- We will consider two views of optimization.
 1. Fixed-point
 2. Updating algorithm

Fixed-point View of Mean Field

- Constrained optimization problem:

$$\begin{aligned} & \max_{Q \in \mathcal{Q}} H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi] \\ & \text{such that } \forall \mathbf{x}, \quad Q(\mathbf{x}) = \prod_i Q_i(x_i) \\ & \quad \forall i \quad \sum_{x_i} Q_i(x_i) = 1 \end{aligned}$$

- Objective is ideal (energy functional); space of distributions is approximate.
 - “Bigger,” more expressive \mathcal{Q} , better approximation.

Restrict to Q_i

$$\max_{Q \in \mathcal{Q}} H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi]$$

$$\text{such that } \forall \mathbf{x}, \quad Q(\mathbf{x}) = \prod_i Q_i(x_i)$$

$$\forall i \quad \sum_{x_i} Q_i(x_i) = 1$$

$$F_i(Q_i) = H_{Q_i} + \sum_{\phi \in \Phi} \mathbb{E}_{Q_i}[\log \phi]$$

Restrict to Q_i

$$\max_{Q \in \mathcal{Q}} H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi]$$

$$\text{such that } \forall \mathbf{x}, \quad Q(\mathbf{x}) = \prod_i Q_i(x_i)$$

$$\forall i \quad \sum_{x_i} Q_i(x_i) = 1$$

$$F_i(Q_i) = H_{Q_i} + \sum_{\phi \in \Phi} \mathbb{E}_{Q_i}[\log \phi]$$

concave in Q_i


linear in Q_i

- Concave function; stationary point is a global maximum (given all other components of Q).
- Next: find that global maximum.

Lagrangian

$$\max_{Q \in \mathcal{Q}} H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi]$$

$$\text{such that } \forall \mathbf{x}, \quad Q(\mathbf{x}) = \prod_i Q_i(x_i)$$


$$\forall i \quad \sum_{x_i} Q_i(x_i) = 1$$


$$L_i(Q_i, \lambda) = H_{Q_i} + \sum_{\phi \in \Phi} \mathbb{E}_{Q_i}[\log \phi] + \lambda \left(\sum_{x_i} Q_i(x_i) - 1 \right)$$

Differentiate for $Q_i(x_i)$

$$\max_{Q \in \mathcal{Q}} H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi]$$

$$\text{such that } \forall \mathbf{x}, \quad Q(\mathbf{x}) = \prod_i Q_i(x_i)$$

$$\forall i \quad \sum_{x_i} Q_i(x_i) = 1$$


$$L_i(Q_i, \lambda) = H_{Q_i} + \sum_{\phi \in \Phi} \mathbb{E}_{Q_i}[\log \phi] + \lambda \left(\sum_{x_i} Q_i(x_i) - 1 \right)$$

$$\frac{\partial L_i}{\partial Q_i(x_i)} = \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi \mid X_i = x_i] - \log Q_i(x_i) - 1 + \lambda$$

$$= 0$$

$$\log Q_i(x_i) = \lambda - 1 + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi \mid X_i = x_i]$$

Stationary Point

$$\max_{Q \in \mathcal{Q}} H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi]$$

$$\text{such that } \forall \mathbf{x}, \quad Q(\mathbf{x}) = \prod_i Q_i(x_i)$$

$$\forall i \quad \sum_{x_i} Q_i(x_i) = 1$$

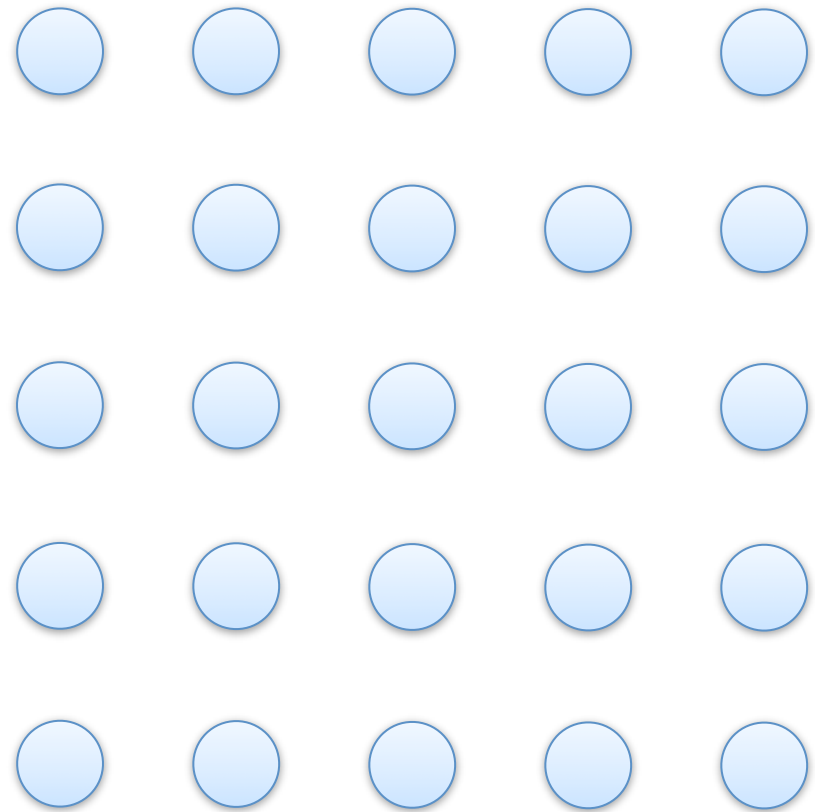
$$Q_i(x_i) = \frac{1}{Z_i} \exp \left(\sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi \mid X_i = x_i] \right)$$

- Gibbs!
- *Does* depend on the other Q_j !
- Fixing other parts of Q , this is a global optimum of F_i .

Example: Mean Field Q

$$Q_i(x_i) = \frac{1}{Z_i} \exp \left(\sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi \mid X_i = x_i] \right)$$

- Each Q_i needs to know about the expected value of P 's log-factors, under different values of X_i .
- Not all factors depend on X_i , of course.



Expectation

$$\max_{Q \in \mathcal{Q}} H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi]$$

$$\text{such that } \forall \mathbf{x}, \quad Q(\mathbf{x}) = \prod_i Q_i(x_i)$$

$$\forall i \quad \sum_{x_i} Q_i(x_i) = 1$$

$$Q_i(x_i) = \frac{1}{Z_i} \exp \left(\sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi \mid X_i = x_i] \right)$$

$$\sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi \mid X_i = x_i] = \sum_{\phi \in \Phi} \mathbb{E}_Q[\log U(X_i, \mathbf{X}_{-i}) \mid X_i = x_i] \quad \backslash \phi(U(X_i, \dots))$$

$$= \sum_{\phi \in \Phi} \mathbb{E}_{Q(\mathbf{x})}[\log U(X_i, \mathbf{X}_{-i}) \mid X_i = x_i]$$

$$= \sum_{\phi \in \Phi} \mathbb{E}_{Q(\mathbf{x}_{-i})}[\log U(x_i, \mathbf{X}_{-i})]$$

$$= \sum_{\phi \in \Phi} \mathbb{E}_{Q(\mathbf{x}_{-i})}[\log ZP(x_i, \mathbf{X}_{-i})]$$

Remove the sum over factors

$$= \sum_{\phi \in \Phi} \mathbb{E}_{Q(\mathbf{x}_{-i})}[\log P(x_i \mid \mathbf{X}_{-i}) + \log ZP(\mathbf{X}_{-i})]$$

$$= \sum_{\phi \in \Phi} \mathbb{E}_{Q(\mathbf{x}_{-i})}[\log P(x_i \mid \mathbf{X}_{-i})] + \mathbb{E}_{Q(\mathbf{x}_{-i})}[\log ZP(\mathbf{X}_{-i})]$$

$$= \sum_{\phi \in \Phi} \mathbb{E}_{Q(\mathbf{x}_{-i})}[\log P(x_i \mid \mathbf{X}_{-i})] + \text{constant}(x_i)$$

Fixed-Point

$$\max_{Q \in \mathcal{Q}} H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi]$$

$$\text{such that } \forall \mathbf{x}, \quad Q(\mathbf{x}) = \prod_i Q_i(x_i)$$

$$\forall i \quad \sum_{x_i} Q_i(x_i) = 1$$

$$Q_i(x_i) = \frac{1}{Z_i} \exp \left(\sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi \mid X_i = x_i] \right)$$

$$Q_i(x_i) = \frac{1}{Z_i} \exp \left(\mathbb{E}_{Q(\mathbf{x}_{-i})}[\log P(x_i \mid \mathbf{X}_{-i})] \right)$$

$$= \frac{1}{Z_i} \exp \left(\sum_{\mathbf{x}_{-i}} \left(\prod_{j \neq i} Q_j(x_j) \right) \log P(x_i \mid \mathbf{x}_{-i}) \right)$$

$$= \frac{1}{Z_i} \prod_{\mathbf{x}_{-i}} \exp \left(\left(\prod_{j \neq i} Q_j(x_j) \right) \log P(x_i \mid \mathbf{x}_{-i}) \right)$$

only worry about
conditional
probability of X_i given
the rest

geometric average of x_i 's conditional probability;

compare to marginal under P , $P(x_i)$, which is an arithmetic average

(Geometric Average)

$$\begin{aligned} \left(\prod_{i=1}^n a_i \right)^{\frac{1}{n}} &= \sqrt[n]{a_1 a_2 \cdots a_n} \\ &= \exp \left(\frac{1}{n} \sum_{i=1}^n \log a_i \right) \end{aligned}$$

From Fixed-point to an Algorithm

- Fixed-point condition only tells us a necessary condition for local optimality.
- It doesn't tell us how to get there.
- Let's get practical!

Only Some Factors Matter

$$\begin{aligned} Q_i(x_i) &= \frac{1}{Z_i} \exp \left(\sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi \mid X_i = x_i] \right) \\ &= \frac{1}{Z_i} \exp \left(\sum_{\phi: X_i \in \text{Scope}(\phi)} \mathbb{E}_{Q(\text{Scope}(\phi) \setminus \{X_i\})}[\log \phi \mid X_i = x_i] \right) \end{aligned}$$

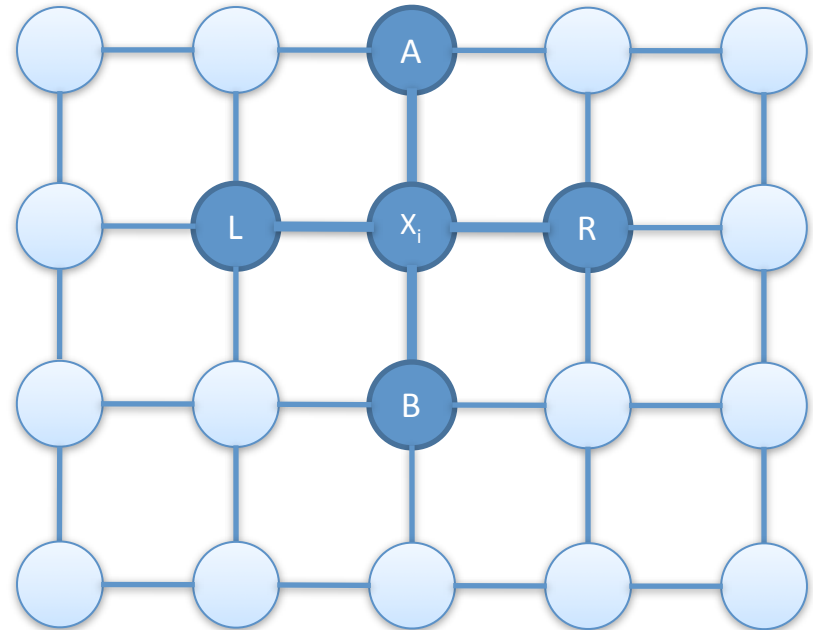
- Q_i needs to be consistent with the expectations of (only!) the potentials in which X_i appears.
 - This is why it's efficient!

Example: One Q_i

$$\phi_i(\text{fg}) = \exp \frac{-\|c_i - \mu_{\text{fg}}\|^2}{\sigma^2}$$

$$\phi_i(\text{bg}) = \exp \frac{-\|c_i - \mu_{\text{bg}}\|^2}{\sigma^2}$$

$$\phi_{i,j}(X_i, X_j) = \begin{cases} 10 & \text{if } X_i = X_j \\ 1 & \text{otherwise} \end{cases}$$



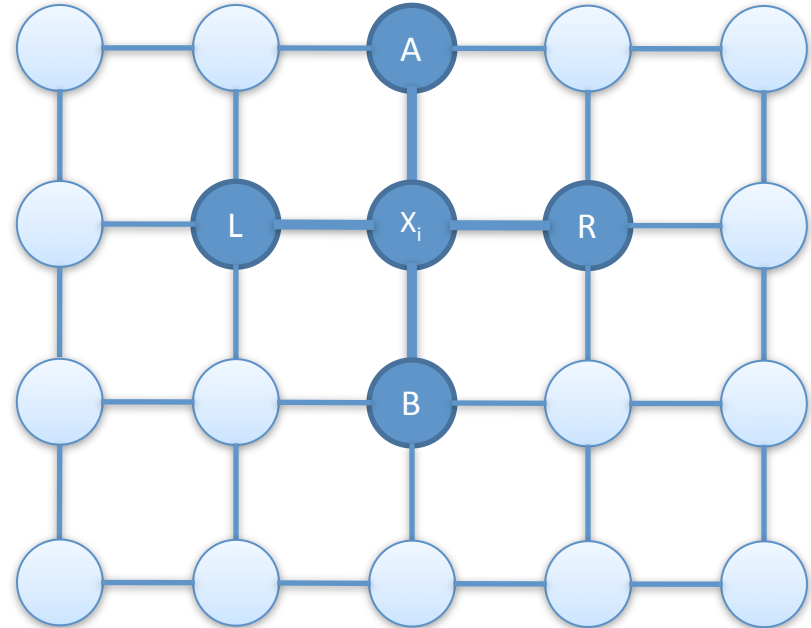
$$Q_{X_i}(\text{fg}) = \frac{1}{Z_i} \exp \left(\begin{array}{l} \log \phi_i(\text{fg}) + \\ Q_A(\text{fg}) \log \phi_{A,X_i}(\text{fg}, \text{fg}) + Q_A(\text{bg}) \log \phi_{A,X_i}(\text{bg}, \text{fg}) + \\ Q_B(\text{fg}) \log \phi_{B,X_i}(\text{fg}, \text{bg}) + Q_B(\text{bg}) \log \phi_{B,X_i}(\text{bg}, \text{fg}) + \\ Q_L(\text{fg}) \log \phi_{L,X_i}(\text{fg}, \text{fg}) + Q_L(\text{bg}) \log \phi_{L,X_i}(\text{bg}, \text{fg}) + \\ Q_R(\text{fg}) \log \phi_{R,X_i}(\text{fg}, \text{fg}) + Q_R(\text{bg}) \log \phi_{R,X_i}(\text{bg}, \text{fg}) \end{array} \right)$$

Example: One Q_i

$$\phi_i(\text{fg}) = \exp \frac{-\|c_i - \mu_{\text{fg}}\|^2}{\sigma^2}$$

$$\phi_i(\text{bg}) = \exp \frac{-\|c_i - \mu_{\text{bg}}\|^2}{\sigma^2}$$

$$\phi_{i,j}(X_i, X_j) = \begin{cases} 10 & \text{if } X_i = X_j \\ 1 & \text{otherwise} \end{cases}$$



$$Q_{X_i}(\text{fg}) = \frac{1}{Z_i} \exp \left(\begin{array}{l} \log \phi_i(\text{fg}) + \\ Q_A(\text{fg}) \log \phi_{A,X_i}(\text{fg}, \text{fg}) + Q_A(\text{bg}) \log \phi_{A,X_i}(\text{bg}, \text{fg}) + \\ Q_B(\text{fg}) \log \phi_{B,X_i}(\text{fg}, \text{bg}) + Q_B(\text{bg}) \log \phi_{B,X_i}(\text{bg}, \text{fg}) + \\ Q_L(\text{fg}) \log \phi_{L,X_i}(\text{fg}, \text{fg}) + Q_L(\text{bg}) \log \phi_{L,X_i}(\text{bg}, \text{fg}) + \\ \boxed{Q_R(\text{fg}) \log \phi_{R,X_i}(\text{fg}, \text{fg}) + Q_R(\text{bg}) \log \phi_{R,X_i}(\text{bg}, \text{fg})} \end{array} \right)$$

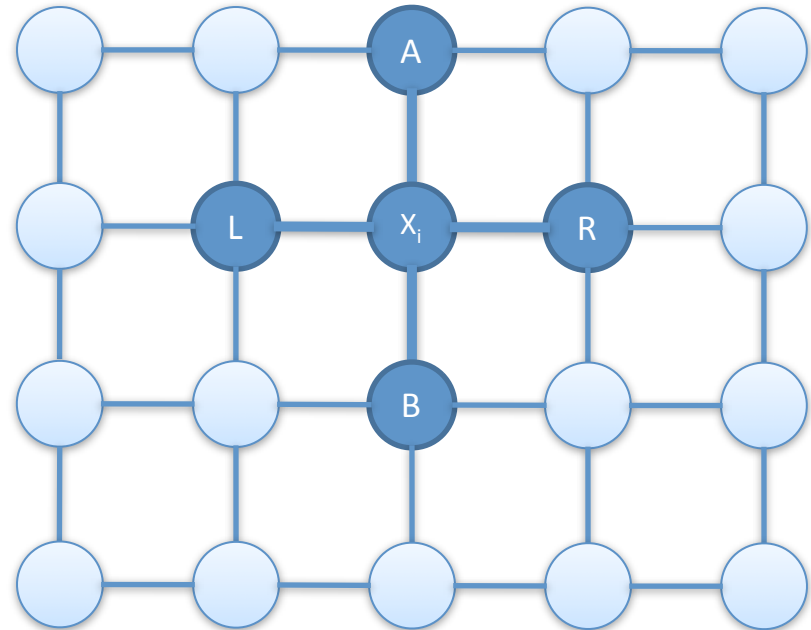
geometric average of the potential between X_i and R

Example: One Q_i

$$\phi_i(\text{fg}) = \exp \frac{-\|c_i - \mu_{\text{fg}}\|^2}{\sigma^2}$$

$$\phi_i(\text{bg}) = \exp \frac{-\|c_i - \mu_{\text{bg}}\|^2}{\sigma^2}$$

$$\phi_{i,j}(X_i, X_j) = \begin{cases} 10 & \text{if } X_i = X_j \\ 1 & \text{otherwise} \end{cases}$$



$$Q_{X_i}(\text{fg}) = \frac{1}{Z_i} \exp \left(\begin{array}{l} \log \phi_i(\text{fg}) + \\ Q_A(\text{fg}) \log \phi_{A,X_i}(\text{fg}, \text{fg}) + Q_A(\text{bg}) \log \phi_{A,X_i}(\text{bg}, \text{fg}) + \\ Q_B(\text{fg}) \log \phi_{B,X_i}(\text{fg}, \text{bg}) + Q_B(\text{bg}) \log \phi_{B,X_i}(\text{bg}, \text{fg}) + \\ Q_L(\text{fg}) \log \phi_{L,X_i}(\text{fg}, \text{fg}) + Q_L(\text{bg}) \log \phi_{L,X_i}(\text{bg}, \text{fg}) + \\ Q_R(\text{fg}) \log \phi_{R,X_i}(\text{fg}, \text{fg}) + Q_R(\text{bg}) \log \phi_{R,X_i}(\text{bg}, \text{fg}) \end{array} \right)$$

no mention of Q_{X_i} !

Inner Step of the Algorithm

- Assuming all other Q_j are fixed, recalculate Q_i :

$$Q_i(x_i) \leftarrow \exp \left(\sum_{\phi: X_i \in \text{Scope}(\phi)} \mathbb{E}_{Q(\text{Scope}(\phi) \setminus \{X_i\})} [\log \phi \mid X_i = x_i] \right)$$

Then renormalize so Q_i sums to one.

- This is a coordinate ascent step on the energy functional.

Mean Field Algorithm

Inputs: Φ and initial value of Q

Output: Q

- Let $\mathbf{U} = \mathbf{X}$
- while \mathbf{U} is not empty:
 - Choose X_i from \mathbf{U} ; store Q_i
 - $Q_i(x_i) \leftarrow \exp \left(\sum_{\phi: X_i \in \text{Scope}(\phi)} \mathbb{E}_{Q(\text{Scope}(\phi) \setminus \{X_i\})} [\log \phi \mid X_i = x_i] \right)$
 - Normalize Q_i
 - If Q_i actually changed, add to \mathbf{U} all variables that share any factor with X_i
 - Remove X_i from \mathbf{U}
- Return Q

Claims

- Every step of the mean field algorithm will improve the energy functional.
- At convergence, we have a stationary point.
 - Could be local minimum, local maximum, or saddle point.
 - In practice, usually a local maximum.

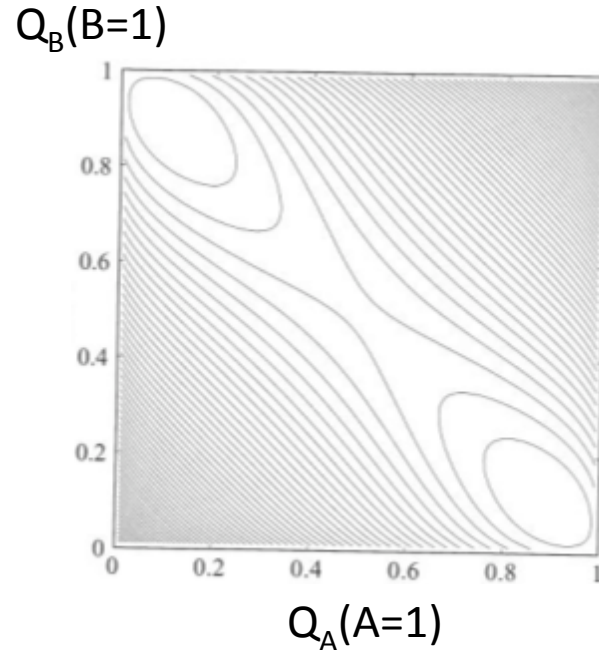
$$\max_{Q \in \mathcal{Q}} H_Q + \sum_{\phi \in \Phi} \mathbb{E}_Q[\log \phi]$$

$$\text{such that } \forall \mathbf{x}, \quad Q(\mathbf{x}) = \prod_i Q_i(x_i)$$

$$\forall i \quad \sum_{x_i} Q_i(x_i) = 1$$

Local Maxima

- This can't be represented by $Q_A(A) \cdot Q_B(B)$
- With small ε , there are two local optima of the energy functional.

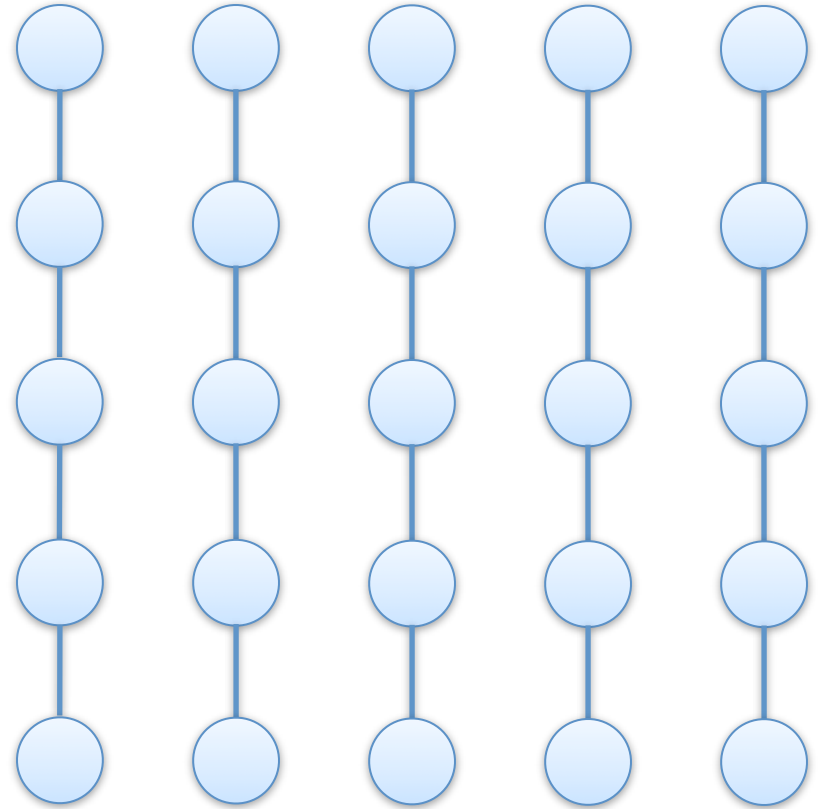


A	B	$\phi(A, B)$
0	0	ε
0	1	$1 - \varepsilon$
1	0	$1 - \varepsilon$
1	1	ε

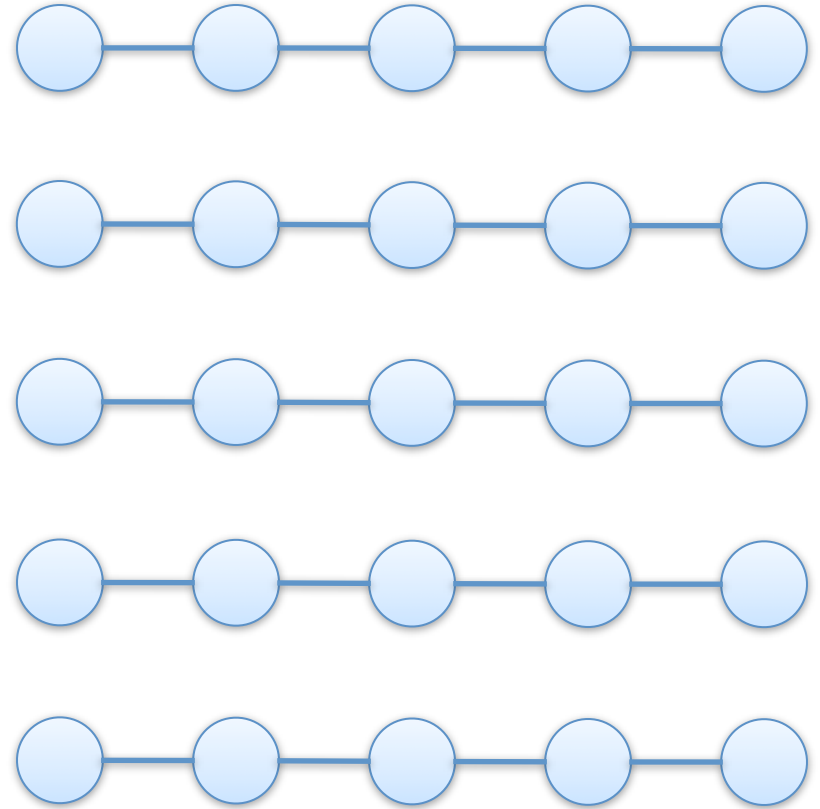
Other Structured Approximations

- Letting \mathcal{Q} be defined by the completely disconnected graph is extreme.
 - Poor approximations.
 - Add some edges?

Example



Example



Other Structured Approximations

- Letting \mathcal{Q} be defined by the completely disconnected graph is extreme.
 - Poor approximations.
 - Add some edges?
- The result of this is a more complex update rule.

$$Q(\mathbf{X}) = \frac{1}{Z_Q} \prod_j \psi_j$$

$$\psi_j(\mathbf{c}_j) \leftarrow \frac{1}{Z_{\psi_j}} \exp \left(\mathbb{E}_Q[\log U \mid \mathbf{c}_j] - \sum_{k \neq j} \mathbb{E}_Q[\log \psi_k \mid \mathbf{c}_j] \right)$$