

# Variational Inference

Sargur Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Plan of discussion

- We first describe inference with PGMs and the intractability of exact inference
- Then give a taxonomy of inference algorithms
- Then the disadvantages of sampling-based inference
- Finally variational inference where we use a tractable distribution that is most similar to the given distribution

# Topics

## 1. Overview

1. Inference
2. Intractability of Inference
3. Shortcomings of sampling-based inference

## 2. Inference as optimization

1. The K-L divergence
2. The variational lower bound
3. On the choice of K-L divergence

## 3. Mean-field inference

# Inference

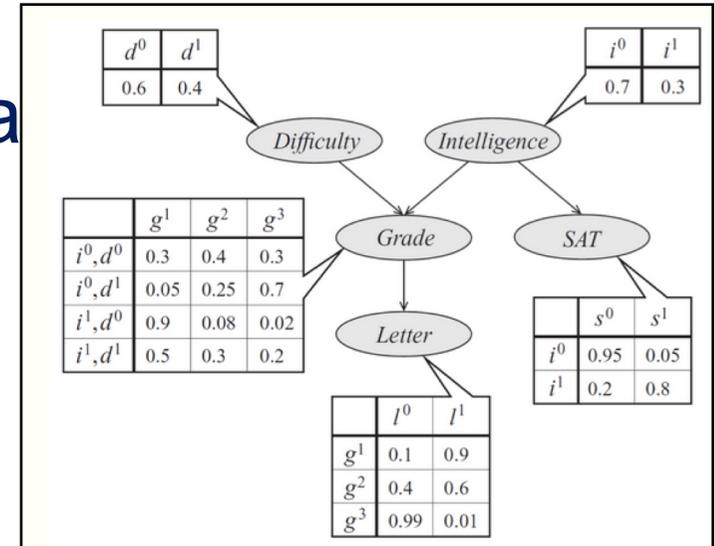
- Knowing *Letter* ( $l=l^0$ ) we have a need to infer *Intelligence* ( $i$ )
  - Can be done using Bayes' rule

$$p(i/l) = p(i, l) / p(l)$$

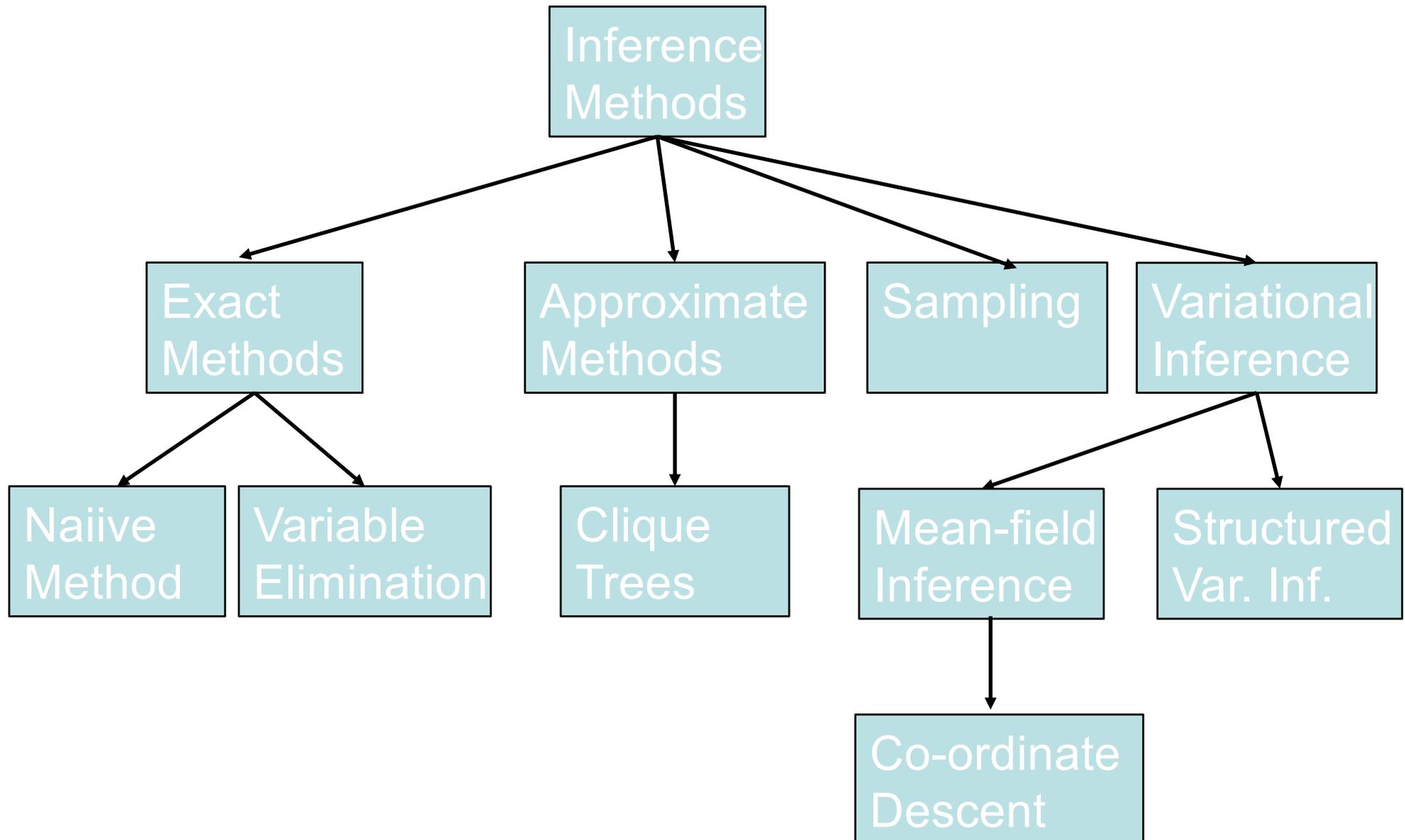
- Marginal  $p(l)$  can be obtained using sum rule

$$p(l) = \sum_g \sum_s \sum_i \sum_d p(l | g) p(s | i) p(i) p(g | i, d) p(d)$$

- With  $n$  variables and  $k$  values each:
  - inference is  $O(k^n)$
  - An intractable problem



# Taxonomy of Inference Algorithms



# Exact Method: Variable Elimination

## 1. Naive Inference

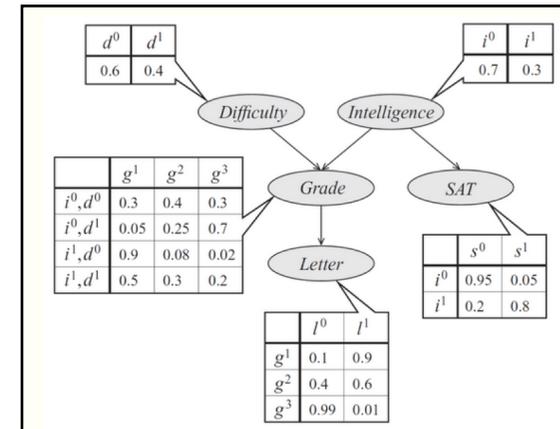
– With  $n$  vars with  $k$  values each, inference is  $O(k^n)$

– Ex: 
$$p(l) = \sum_g \sum_s \sum_i \sum_d p(l | g)p(s | i)p(i)p(g | i, d)p(d)$$

## 2. Variable elimination (VE)

–  $O(nk^{M+1})$

- $M$  is max size of factor  $\tau$
- formed during elimination process



$$p(l) = \sum_g p(l | g) \sum_s \sum_i p(s | i)p(i) \underbrace{\sum_d p(g | i, d)p(d)}_{\tau_1(g, i)} \underbrace{\phantom{\sum_d p(g | i, d)p(d)}}_{\tau_2(g, s)} \underbrace{\phantom{\sum_d p(g | i, d)p(d)}}_{\tau_3(g)}$$

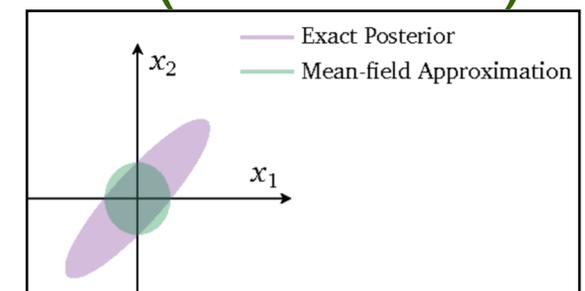
Choosing optimal VE Ordering is NP-hard  
 Heuristics often used:  
 Min-nbrs, Min wt, Min-fill

# Shortcoming of Sampling inference

- Most sampling-based inference algorithms are instances of MCMC, e.g.,
  - Gibbs sampling and Metropolis-Hastings
- Important shortcomings
  - Although guaranteed to find global optimum, cannot tell how close to a good solution in finite time
  - To reach good solution, MCMC methods require choosing an appropriate sampling technique, e.g., a good proposal in MCMC
    - Choosing the sampling technique is an art in itself

# Variational inference: basic idea

- An alternative approach to inference:
  - Called the variational family of algorithms
  - Cast inference as an optimization problem
- We are given an intractable distribution  $p$ 
  1. Find a tractable distribution  $q \in Q$  most similar to  $p$ 
    - Ex: Tractable distribution:  $q(\mathbf{x}) = q_1(x_1)q_2(x_2)\dots q_n(x_n)$   
where  $q_i(x_i)$  is a 1-dimensional discrete distribution
    - Similarity measured by  $KL(p||q)$
  2. We query  $q$  to get approximate solution (inference)



# Sampling vs variational

## 1. Convergence

- Unlike sampling, variational inference will almost never find the optimum solution
- However we always know whether converged
- We may even have bounds on accuracy

## 2. Variational inference advantages

1. Scale better
  2. Amenable to SGD
  3. Parallelization over multiple processors
  4. Acceleration using GPUs
- Variational is now more popular than sampling

# The Kullback-Leibler divergence

- K-L divergence between distributions  $p$  and  $q$  is given by

$$KL(q||p) = \sum_x q(x) \log \frac{q(x)}{p(x)}$$

- In information theory,  $KL$  measures differences in information contained within two distributions
- Properties that make it especially useful here:
  - $KL(q||p) \geq 0$  for all  $q, p$
  - $KL(q||p) = 0$  if and only if  $q = p$
  - Since  $KL(q||p) \neq KL(p||q)$  it is divergence, not distance

# Variational Lower Bound for $Z(\theta)$

- Suppose distribution  $p$  over  $\mathbf{x} = [x_1, \dots, x_n]$  is

$$p(\mathbf{x}; \theta) = \frac{\tilde{p}(\mathbf{x}; \theta)}{Z(\theta)} = \frac{1}{Z(\theta)} \prod_{k \in C} \phi_k(x_k; \theta)$$

$$Z(\theta) = \sum_{x_1, \dots, x_n} \prod_{k \in C} \phi_k(x_k; \theta)$$

- It is intractable because of  $Z(\theta)$

- Consider  $J(q) = KL(q \parallel \tilde{p}) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{\tilde{p}(\mathbf{x})}$  which is tractable

– Its negation is the lower bound for  $Z(\theta)$  as follows

$$\begin{aligned}
 J(q) = KL(q \parallel \tilde{p}) &= \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{\tilde{p}(\mathbf{x})} \\
 &= \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} - \log Z(\theta) \\
 &= KL(q \parallel p) - \log Z(\theta)
 \end{aligned}
 \Rightarrow \log Z(\theta) = KL(q \parallel p) - J(q) \geq -J(q)$$

Note:

$$\tilde{p}(x) = p(x)Z(\theta)$$

- i.e., a tractable  $-J(q)$  is lower bound for  $\log Z(\theta)$

# Summary of Optimization

- We want to maximize  $-J(q)$
- Thus the optimization problem is

$$\min_q J(q)$$

$$J(q) = \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)}$$

- By minimizing  $J(q) = KL(q||p) - \log Z(\theta)$  we are squeezing the difference between  $KL(q||p)$  and its lower bound  $\log Z(\theta)$ 
  - where  $q$  is chosen such that  $J(q)$  is tractable

# Evidence Lower Bound

- In the context of inferring a marginal probability  $p(\mathbf{x}|D) = p(\mathbf{x}, D) / p(D)$ , i.e., of variables  $\mathbf{x}$  given data  $D$ , we can show that
  - Minimizing  $J(q)$  implies maximizing the log-likelihood  $\log p(D)$  of the observed data (evidence)
- Which means that we should choose that model  $q(\mathbf{x})$  which assigns the highest probability to the data  $D$
- Shown next

# Evidence Lower Bound (ELBO)

- Suppose we have joint distribution  $p(X, Z)$

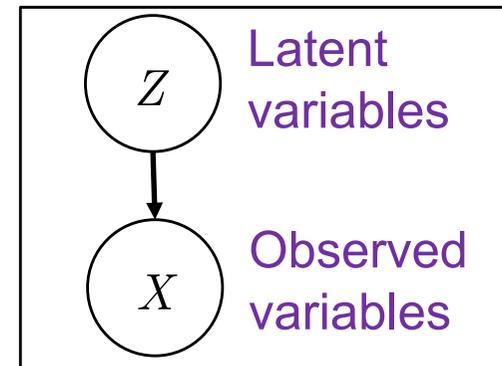
– We are interested in computing

- By approximating  $p(Z|X)$  by  $q(Z)$

$$\begin{aligned}
 KL[q(Z)||p(Z|X)] &= \int_Z q(Z) \log \frac{q(Z)}{p(Z|X)} \\
 &= - \int_Z q(Z) \log \frac{p(Z|X)}{q(Z)} \\
 &= - \left( \int_Z q(Z) \log \frac{p(X, Z)}{q(Z)} - \int_Z q(Z) \log p(X) \right) \\
 &= - \int_Z q(Z) \log \frac{p(X, Z)}{q(Z)} + \log p(X) \int_Z q(Z) \\
 &= -L + \log p(X) \quad \text{where lower bound}
 \end{aligned}$$

Rearranging  $L = \log p(X) - KL[q(Z)||p(Z|X)]$

$$p(Z|X) = \frac{p(X|Z)p(Z)}{p(X)} = \frac{p(X|Z)p(Z)}{\int_Z p(X, Z)}$$



$$L = \mathbb{E}_q[\log p(X, Z)] + H[Z]$$

– Minimizing  $J(q) = -L$  implies maximizing  $\log p(X)$  of observed

- Hence  $-J(q)$  is known as the evidence lower bound (ELBO)

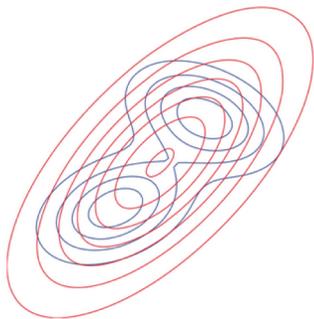
– commonly written as  $\log Z(\theta) \geq \mathbb{E}_{q(x)}[\log \tilde{p}(x) - \log q(x)]$

- Difference between  $\log Z(\theta)$  and  $-J(q)$  is  $KL(q||p)$

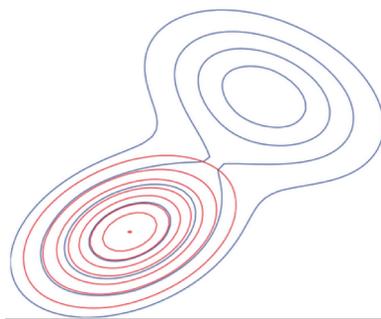
- Thus by maximizing ELBO we are minimizing  $KL(q||p)$

# On the choice of KL divergence

- For variational inference
  - Maximizing lower bound leads to minimize  $KL(q||p)$
- Since  $KL(q||p) \neq KL(p||q)$  how to choose one
  - Difference is mostly computational
  - Involves expectation wrt  $q$
  - Involves expectation wrt  $p$  which is intractable

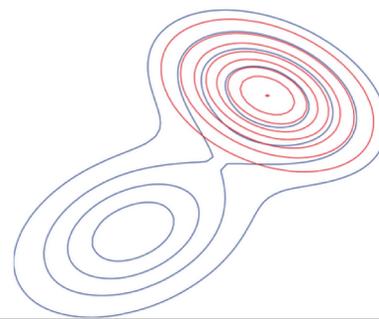


Fitting a unimodal  $q$  (red) to a multimodal  $p$  (blue)



$KL(p||q)$  tries to cover both modes

Inclusive  
K-L Divergence



$KL(q||p)$  forces  $q$  to choose one of the modes of  $p$

Exclusive  
K-L Divergence

# Mean Field Inference

- For choice of the approximating family  $Q$ 
  - ML literature contains dozens of methods for parameterizing distributions:
    - exponential families, neural nets, Gaussian processes, etc
- Most widely used class is the fully factored set
  - $q(\mathbf{x}) = q_1(x_1)q_2(x_2)\dots q_n(x_n)$  where  $q_i(x_i)$  is a distribution over a 1-D discrete variable
- The optimization problem is

$$\min_{q_1, \dots, q_n} J(q)$$

Note that  $J(q) = KL(p||q) - \log Z(\theta)$  and

$$p(x; \theta) = \frac{\tilde{p}(x; \theta)}{Z(\theta)} = \frac{1}{Z(\theta)} \prod_k \phi_k(x_k, \theta_k)$$

$$J(q) = \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)}$$

is tractable

# Popularity of mean field inference

- The choice of  $q(\mathbf{x}) = q_1(x_1) q_2(x_2) \dots q_n(x_n)$  works surprisingly well
- It is the most popular choice of optimizing the variational bound
- Variational inference with this choice is known as mean field inference

# Performing mean field optimization

- Done via coordinate descent over  $q_i$
- We iterate over  $j=1,2,\dots,n$ 
  - For each  $j$  we optimize  $KL(q||p)$  over  $q_j$  while keeping the other co-ordinates  $q_{-j} = \prod_{i \neq j} q_i$  fixed
  - The optimization problem for one coordinate has a closed-form solution

$$\log q_j(x_j) \leftarrow \mathbb{E}_{q_{-j}} [\log \tilde{p}(x)] + \text{const.}$$

- Both sides of equation contain univariate functions of  $x_j$
- We are thus replacing  $q(x_j)$  with another function of same form
- The constant term is a normalization constant for the new distribution

# What is Coordinate descent?

- Convex optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

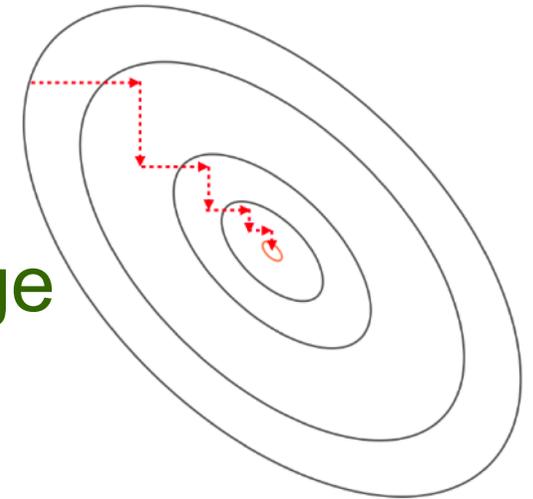
- Where  $f$  is differentiable and  $n$  is large

- Coordinate descent:

- Select a coordinate to update
- Take a small gradient step along coordinate

- Coordinate descent is  $n$  times faster than gradient descent when  $f$  has the form

$$g(x) = \underbrace{\sum_{i=1}^n f_i(x_i)}_{\text{separable}} + \underbrace{\sum_{i=1}^n \sum_{j=1}^n f_{ij}(x_i, x_j)}_{\text{pairwise separable}} + \underbrace{f(Ax)}_{\text{linear composition}}$$



# Mean Field Optimization

- The optimization problem is  $\min_{q_1, \dots, q_n} J(q)$
- Optimization for one coordinate has the simple closed form solution

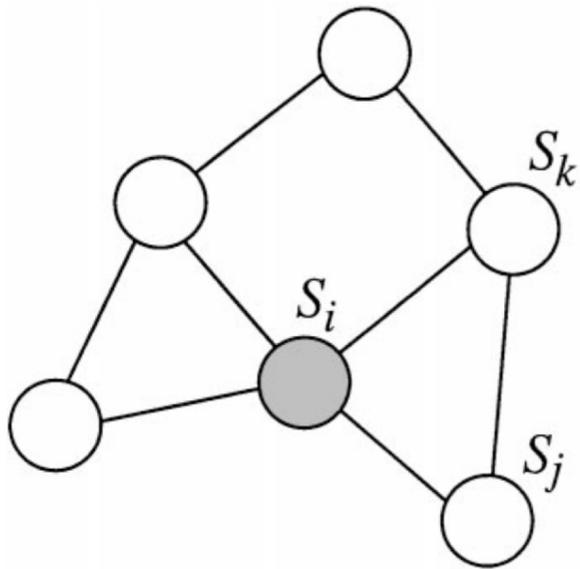
$$\log q_j(x_j) \leftarrow \mathbb{E}_{q_{-j}} [\log \tilde{p}(x)] + \text{const}$$

- The rhs is taking an expectation of a sum of factors

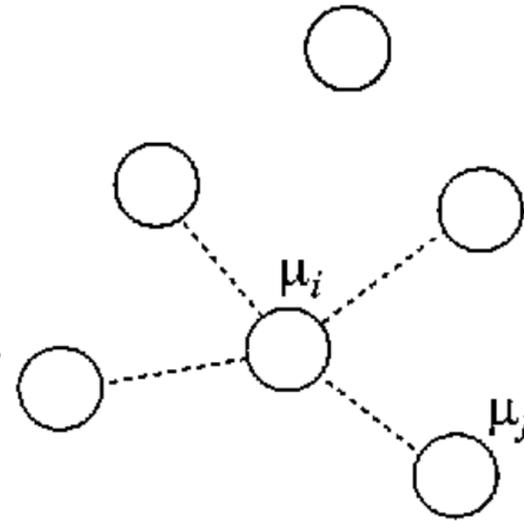
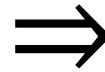
$$\log \tilde{p}(x) = \sum_k \log \phi(x_k)$$

- Of these, only factors belonging to the Markov blanket of  $x_j$  are a function of  $x_j$  (by definition of Markov blanket)
- the rest are constant wrt  $x_j$  and can be pushed into the constant term.

# Markov Blanket & Mean Field



A node  $S_i$  in a Boltzmann Machine and its Markov blanket



Approximating mean field distribution  $q$ :  
 based on a graph with no edges.  
 Mean field yields a deterministic relationship  
 between the variational parameters  
 $\mu_i$  and  $\mu_j$  for nodes  $j$  in the Markov blanket of  $S_i$ .

# Complexity of Mean Field Inference

- If variables are discrete with  $K$  possible values,
- There are  $F$  factors and  $N$  variables in the Markov blanket of  $x_j$ ,
- Then computing expectation is  $O(KFK^N)$  :
  - for each value of  $x_j$  we sum over all  $K^N$  assignments of the  $N$  variables, and in each case, we sum over the  $F$  factors.

# Result of procedure

- The result iteratively fits a fully-factored  $q(\mathbf{x}) = q_1(x_1) q_2(x_2) \cdots q_n(x_n)$  that approximates  $p$  in terms of  $KL(q\|p)$
- After each step of coordinate descent, we increase the variational lower bound, tightening it around  $\log Z(\theta)$

# Structured Variational Approximation

- Main parameter in maximization problem is choice of family  $Q$ 
  - This choice introduces a trade-off
    - Families that are simpler, i.e., BNs and MNs of small tree width allow more efficient inference
    - If  $Q$  is too restrictive then it cannot represent distributions that are good approximations of  $P_{\Phi}$ 
      - Giving rise to a poor approximation  $Q$
- Family is chosen to have enough structure
  - Hence called *structured variational approximation*
    - Variational calculus since we maximize over functions
    - Unlike belief propagation guaranteed to lower bound the log-partition function and guaranteed to converge

# References

- D. Koller and N. Friedman, *Probabilistic Graphical Models*, MIT Press, 2009
- <https://ermongroup.github.io/cs228-notes/inference/variational/>
- M. I. Jordan, Z. Ghahramani, L. K. Saul, An Introduction to Variational Methods for Graphical Models, *Machine Learning* 1999