

# Mean Field / Variational Approximations

Presented by Jose Nuñez

10/24/05

## Outline

- Introduction
- Mean Field Approximation
- Structured Mean Field
- Weighted Mean Field
- Variational Methods

## Introduction

Problem:

- We have distribution  $P(x)$  but inference is hard to compute.

Previous solutions:

- Approximate energy functional: Bethe, Kikuchi

## Introduction

New idea:

- Directly optimize the energy functional introducing a distribution  $Q(x)$  defined on the same domain of variables as  $P$  which incorporates some constraints.
- **Objective:** We want to find  $Q(x)$  which is the best approximation of  $P(x)$  and use  $Q(x)$  to make inferences.
- Find  $Q \in \mathcal{Q}$  that minimizes  $\mathbf{F}(\mathbf{P}, \mathbf{Q})$

# Mean Field Approximation

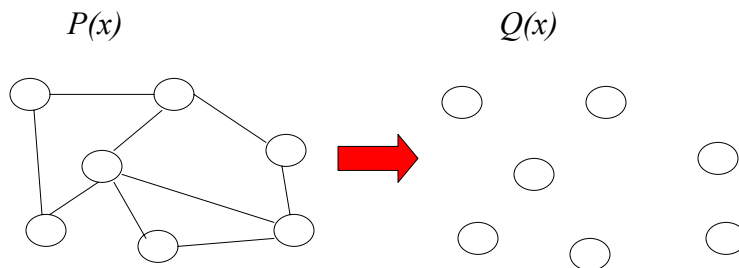
Assumptions:

- $Q(x)$  is our mean field approximation.
- Variables in the  $Q$  distribution are independent variables  $X_i$ .
- In the standard mean field approach,  $Q$  is completely factorized:

$$Q(x) = \prod_i Q_i(x_i)$$

# Mean Field Approximation

What happens when we apply mean field?



## Mean Field Approximation

$$F(P, Q) = -\sum_{\phi \in F} \sum_{x_\phi} Q(x_\phi) \log \phi(x_\phi) + \sum_x Q(x) \log Q(x)$$

$$Q(x) = \prod_i Q_i(x_i)$$

$$E(P, Q) = -\sum_{\phi \in F} \sum_{x_\phi} Q(x_\phi) \log \phi(x_\phi) = -\sum_{\phi \in F} \sum_{x_\phi} \left( \prod_{x_i \in x_\phi} Q(x_i) \right) \log \phi(x_\phi)$$

$$\begin{aligned} H(Q) &= -\sum_x Q(x) \log Q(x) = -\sum_x \left( \prod_{i \in x} Q(x_i) \right) \log \left( \prod_{i \in x} Q(x_i) \right) \\ &= -\sum_x \left( \prod_i Q(x_i) \right) \sum_i \log Q(x_i) \\ &= -\sum_i \sum_{x_i} Q(x_i) \log Q(x_i) \\ &= \sum_i H_{Q_i}(x_i) \end{aligned}$$

## Mean Field Approximation

$$F(P, Q) = -\sum_{\phi \in F} \sum_{x_\phi} Q(x_\phi) \log \phi(x_\phi) + \sum_x Q(x) \log Q(x)$$

$$E(P, Q) = -\sum_{\phi \in F} \sum_{x_\phi} \left( \prod_{x_i \in x_\phi} Q(x_i) \right) \log \phi(x_\phi)$$

$$H(Q) = -\sum_i \sum_{x_i} Q(x_i) \log Q(x_i)$$

**Task: find**  $Q(x) = \prod_i Q_i(x_i)$  **minimizing**  $F(P, Q)$

**such that**  $\sum_{x_i} Q(x_i) = 1$

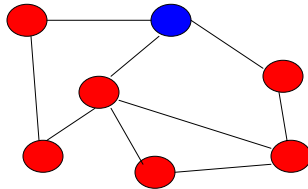
**Solving:** build a Lagrangian, differentiate and set to 0 !

# Mean Field Approximation

The distribution  $Q(x_i)$  is locally optimal solution given  $Q(x_1), \dots, Q(x_{i-1}), Q(x_{i+1}), \dots, Q(x_n)$ , if:

$$Q(x_i) = \frac{1}{Z_i} \exp \left\{ \sum_{\phi \in F} E_Q [\ln \phi | x_i] \right\} \quad \text{MF-equation}$$

Where  $Z_i$  is a local normalizing constant and  $E_Q[\ln \phi | x_i]$  is the conditional expectation given the value  $x_i$ .



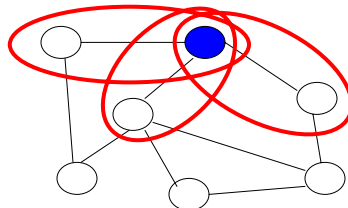
# Mean Field Approximation

Locality:

- Only local operations are needed for iteration of the MF-equations.
- In other words, only neighboring variables are needed.

$$Q(x_i) = \frac{1}{Z_i} \exp \left\{ \sum_{\phi: X_i \in \text{Scope}[\phi]} E_Q [\ln \phi(U_\phi, x_i)] \right\} \quad \text{MF-equation simplified}$$

where  $U_\phi = \text{Scope}[\phi]$



Calculation of  $Q(x_i)$  depends only on clusters  $X_i$  belongs to

# Mean Field Approximation

**Solution:** Iterate mean field equations

$$Q(x_i) = \frac{1}{Z_i} \exp \left\{ \sum_{\phi: X_i \in \text{Scope}[\phi]} E_Q [\ln \phi(U_\phi, x_i)] \right\} \quad \begin{array}{l} \text{MF-equation} \\ \text{simplified} \end{array}$$

- Converge to a fixed point.

**Problem:** convergence to a local optima.

# Mean Field Approximation

Haft et al. paper:

- Optimize the KL divergence instead of the free energy

$$D(Q | P) = E_Q \left( \log \frac{Q(x)}{P(x)} \right)$$

$$D(Q | P) = E_Q (\log Q(x)) - E_Q (\log P(x))$$

$$D(Q | P) = E_Q (\log Q(\bar{X}_i)) - E_Q (\log P(\bar{X}_i)) \\ + E_Q (\log Q(X_i)) - E_Q (\log P(X_i | \bar{X}_i))$$

Assume:  $P(X) = P(X_i | \bar{X}_i)P(\bar{X}_i)$

$$D(Q | P) = E_{Q(\bar{X}_i)} (\log Q(\bar{X}_i)) - E_{Q(\bar{X}_i)} (\log P(\bar{X}_i)) \\ + E_{Q(X_i)} (\log Q(X_i)) - E_{Q(X)} (\log P(X_i | \bar{X}_i))$$

# Mean Field Approximation

Haft et al. paper:

- Optimize the KL divergence instead of the free energy

$$D(Q | P) = E_{Q(\bar{X}_i)}(\log Q(\bar{X}_i)) - E_{Q(\bar{X}_i)}(\log P(\bar{X}_i))$$

$$+ E_{Q(X_i)}(\log Q(X_i)) - E_{Q(X)}(\log P(X_i | \bar{X}_i))$$

Does not depend on  $X_i$

Depends on  $X_i$

$$\frac{\partial D(Q | P)}{\partial Q(X_i)} = \frac{\partial}{\partial Q(X_i)} E_{Q(X_i)}(\log Q(X_i)) - E_{Q(X)}(\log P(X_i | \bar{X}_i))$$

Subject to  $\sum_{X_i} Q(X_i) = 1$

# Mean Field Approximation

Haft et al. paper:

$$Q(X_i) \propto \exp(E_{Q(\bar{X}_i)}(\log P(X_i | \bar{X}_i))) = \exp\langle \log P(X_i | \bar{X}_i) \rangle_{Q(\bar{X}_i)}$$

MF-equation

Locality:

$$Q(x_i) \propto \exp\langle \log P(x_i | M_i) \rangle_{Q(M_i)}$$

MF-equation simplified

where M is the Markov boundary.

# Mean Field Approximation

Algorithm:

```
Procedure Mean-Field (  
   $\mathcal{F}$ , // factors that define  $P_{\mathcal{F}}$   
   $Q_0$  // Initial choice of  $Q$   
)  
   $Q \leftarrow Q_0$   
   $Unprocessed \leftarrow \mathcal{X}$   
  while  $Unprocessed \neq \emptyset$   
    Choose  $X_i$  from  $Unprocessed$   
     $Q_{old}(X_i) \leftarrow Q(X_i)$   
    for  $x_i \in Val(X_i)$  do  
       $Q(x_i) \leftarrow \exp \left\{ \sum_{\phi: X_i \in Scope[\phi]} E_Q[\ln \phi | x_i] \right\}$   
    Normalize  $Q(X_i)$  to sum to one  
    if  $Q_{old}(X_i) \neq Q(X_i)$  then  
       $Unprocessed \leftarrow Unprocessed \cup \left( \bigcup_{\phi: X_i \in Scope[\phi]} Scope[\phi] \right)$   
       $Unprocessed \leftarrow Unprocessed - \{X_i\}$   
  return  $Q$ 
```

# Mean Field Approximation

- Converges to one of typically many local minima.
- Easy to compute but sometimes is not good enough.
- It cannot describe complex posteriors (eg. XOR)
- We must use a richer class of distributions  $Q$ .

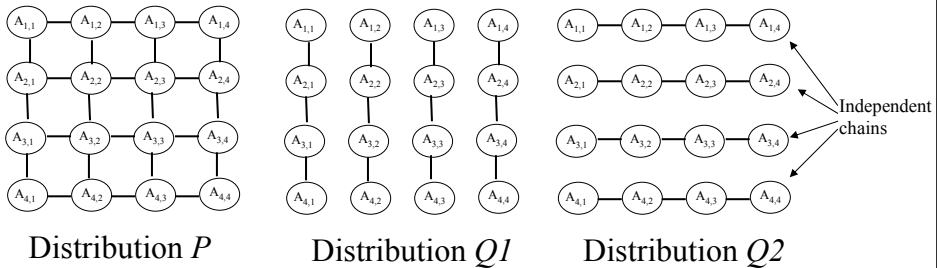


# Structured Mean Field

## Exploiting Substructures

- If we use a distribution  $Q$  that can capture some of the dependencies in  $P$ , we can get a better approximation.

Two possible substructures for  $Q$



# Structured Mean Field

## Exploiting Substructures

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

where  $\psi_j$  is a factor with  $Scope[\psi_j] = C_j$ .

and assume we have the set of potential scopes:

$$\{C_j \subseteq \mathcal{X}: j = 1, \dots, J\}$$

# Structured Mean Field

## Exploiting Substructures

$$\text{Given: } Q(x) = \frac{1}{Z_Q} \prod_j \psi_j$$

$$\text{And restriction: } \sum_{c_j} \psi_j(c_j) = 1$$

Then the potential  $\psi_j$  is locally optimal when:

$$\psi_j(c_j) \propto \exp \left\{ E_Q [\ln P'_r | c_j] - \sum_{k \neq j} E_Q [\ln \psi_k | c_j] \right\}$$

# Structured Mean Field

## Exploiting Substructures

- Locality as Mean Fields:

$$\psi_j(c_j) \propto \exp \left\{ \sum_{\phi \in A_j} E_Q [\phi | c_j] - \sum_{\psi_k \in B_j} E_Q [\ln \psi_k | c_j] \right\}$$

where

$$A_j = \{ \phi \in F : Q \not\perp (U_\phi \perp C_j) \}$$

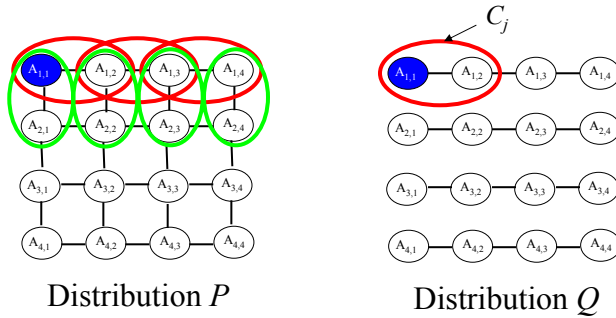
and

$$B_j = \{ \psi_k : Q \not\perp (C_k \perp C_j) \}$$

# Structured Mean Field

Updating:

- Calculation of  $Q(X_i)$  depends on clusters where  $X_i$  belongs to. And on clusters overlapping  $C_j$  (in  $Q$ ). And on scopes  $C_k$  dependent of  $C_j$  (in  $Q$  also).



# Structured Mean Field

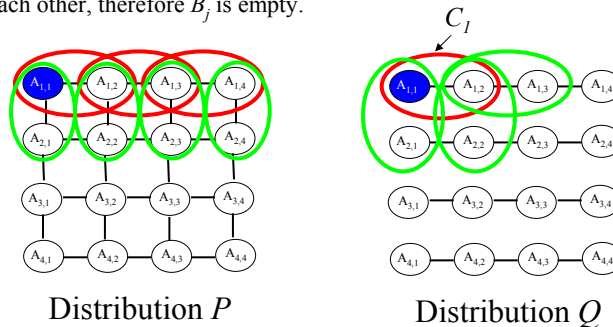
In other words, we want to compute  $A_{1,1}$ :

$$C_1 = \{A_{1,1}, A_{1,2}\} \quad C_2 = \{A_{1,2}, A_{1,3}\} \quad C_3 = \{A_{2,1}, A_{2,2}\} \dots$$

$A_j$  = Clusters  $X_i$  belongs to (as standard mean field) i.e.  $\{A_{1,1}, A_{1,2}\}$  and  $\{A_{1,1}, A_{2,1}\}$

Clusters overlapping  $C_j$  and those from PF. For example in this case  $A_{1,2}$  in  $C_1$  overlaps in  $P$ , thus we need to consider  $\{A_{1,2}, A_{1,3}\}$  and  $\{A_{1,1}, A_{2,2}\}$ . The same occurs with  $A_{1,3}$  and  $A_{1,4}$ .

$B_j$  = Clusters in  $Q$  dependent on  $C_j$ . In this example every  $C$  is independent from each other, therefore  $B_j$  is empty.



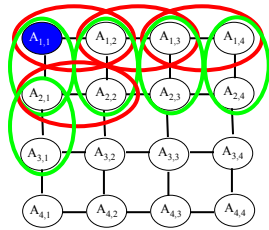
# Structured Mean Field

Again we want to compute  $A_{1,1}$ , assume the new substructure in  $Q$ :

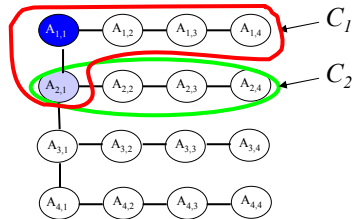
Now we choose  $C_1 = \{A_{1,1} A_{1,2} A_{1,3} A_{1,4} A_{2,1}\}$   $C_2 = \{A_{2,1} A_{2,2} A_{2,3} A_{2,4}\}$

$A_j =$  We consider the same clusters as before but now we add those overlapping with  $A_{2,1}$ , i.e.  $\{A_{2,1} A_{3,1}\}$  and  $\{A_{2,1} A_{2,2}\}$

$B_j =$  Clusters in  $Q$  dependent on  $C_j$ . Now we have  $A_{2,1}$  (in  $C_1$ ) overlapping with  $A_{2,1}$  (in  $C_2$ ). We need to subtract  $A_{2,1}$  since we already used it in  $A_j$ .



Distribution  $P$



Distribution  $Q$

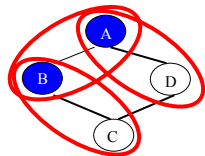
# Structured Mean Field

Another example, we want to compute  $a, b$ :

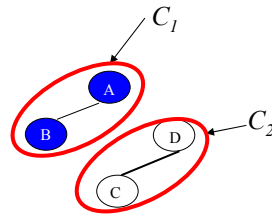
Now we choose  $C_1 = \{A B\}$   $C_2 = \{C D\}$

$A_j = \{\{A B\} \{A D\} \{B C\}\}$

$B_j =$  Empty, since  $C_1$  and  $C_2$  do not overlap.



Distribution  $P$



Distribution  $Q$

# Structured Mean Field

## Exploiting Substructures

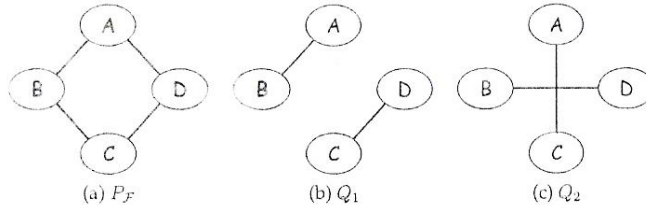
- Updates are relatively costly due to the consideration of structure.

## Two approaches for updates:

- **Sequential:** Choose a factor and update it, then perform inference. It will converge.
- **Parallel:** Update all factors, then inference. It doesn't guarantee convergence.

# Structured Mean Field

Example:



Structure (b) can be exploited:

$$Q(A, B, C, D) = \frac{1}{Z_Q} \psi_1(A, B) \psi_2(C, D)$$

$$Q'(A, B, C, D) = \frac{1}{Z_Q} \phi_{AB}(A, B) \phi_{CD}(C, D) \psi'_1(A) \psi''_1(B) \psi'_2(C) \psi''_2(D)$$

Structure (c) **cannot** be exploited (redundant)

## Structured Mean Field

Refinement Theorem:

- Refines an initial approximating network by factorizing its factors into a product of factors and potentials from  $P_F$ .
- $\psi_k$  can be written as the product of two sets of factors:
  - Those in  $P_F$  that are subsets of the scope of  $\psi_k$ .
  - Partially “covered” factors in  $P_F$  by the scope of  $\psi_j$  and other factors in  $Q$ .

## Weigthed Mean Field

General Mixture Weights

- Idea:

Instead of selecting one particular MF solution, we form a weighted average (a mixture) of several solutions.
- Enumerate the different MF-solutions by a hidden variable  $a$ ,  $Q(X|a)$ .
- Assign mixture weights  $Q(a)$ .

$$Q(X) = \sum_a Q(X|a)Q(a)$$

## Weigthed Mean Field

Given  $Q(X) = \sum_a Q(X|a)Q(a)$

under the constraint  $\sum_a Q(a) = 1$

Determine  $Q(a)$  such that  $D(Q||P)$  is minimized:

$$Q(a) \propto \exp \left[ - \left\langle \log \frac{P(X|a)}{P(X)} \right\rangle_{Q(X|a)} \right]$$
$$\propto \exp [ - D(Q(X|a) || P(X)) ]$$

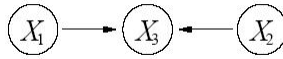
## Weigthed Mean Field

### General Mixture Weights

- The previous formula means that different solutions  $Q(X|a)$  contribute to the global distribution  $Q(X)$  according to their distance to  $P(X)$ .
- Note however, we are **not** optimizing  $Q(X|a)$  simultaneously.

# Weighted Mean Field

Example: Noisy-OR



	first MF-solution $Q(\cdot X_3=1, a=1)$	second MF-solution $Q(\cdot X_3=1, a=2)$	marginals of MF-mixture $Q(\cdot X_3=1)$	exact marginals $P(\cdot X_3=1)$
$X_1=1$	0.137	0.973	0.555	0.528
$X_2=1$	0.973	0.137	0.555	0.528

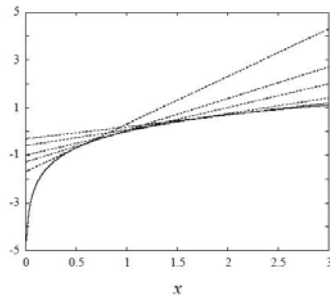
$P(X_1, X_2 X_3=1)$ compared to $Q(X_1, X_2 X_3=1)$	$X_2 = 0$		$X_2 = 1$	
	P:	Q:	P:	Q:
$X_1 = 0$	0.005	0.023	0.466	0.422
$X_1 = 1$	0.466	0.422	0.063	0.133

# Variational Methods

Idea:

- Introducing auxiliary variational parameters that help in simplifying a complex objective function.

$$\ln(x) \leq \lambda x - \ln(\lambda) - 1$$



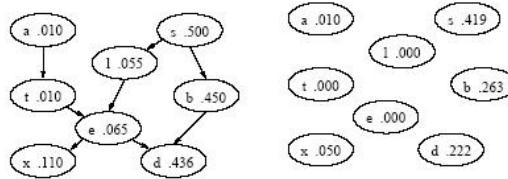
This upper bound allows to approximate  $\ln(x)$  with a term that is linear in  $x$ .



Thank you!

## Mean Field Approximation

Example from Wiegerinck:



(a) Target distribution

(b) KL = 0.43

Noisy-OR from Haft et al.:

	first MF-solution $Q(\cdot X_3=1)$	exact marginals $P(\cdot X_3=1)$
$X_1=1$	0.137	0.528
$X_2=1$	0.973	0.528