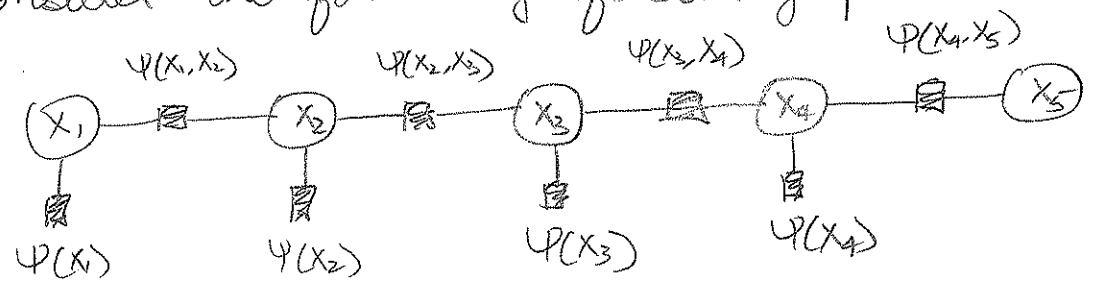


4/1/19

(1)

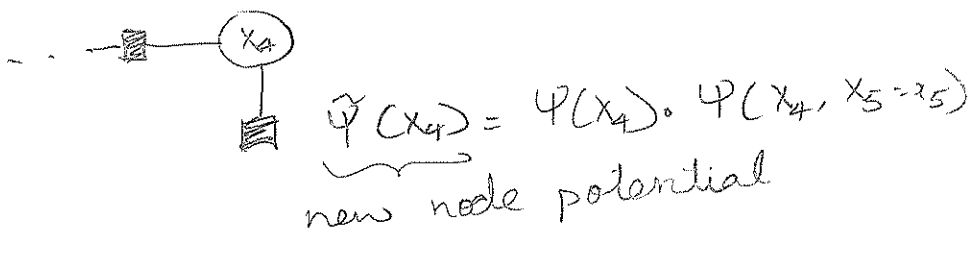
MRF INFERENCE: CONNECTING BP & VE

① Consider the following factor graph:



Query: Find  $P(x_1 | x_5 = z_5) \propto P(x_1, x_5 = z_5)$

Step 1: Instantiate Evidence:



Step 2: Running VE  $\rightarrow$  with Elim Ordering  $\{x_4, x_3, x_2\}$

& BP  $\rightarrow$  with message schedule  $\{(4,3), (3,2), (2,1), (1,2), (2,3), (3,4)\}$

Iter	VE	BP
1	$\sum_{x_4} \psi(x_4) \psi(x_3, x_4)$	$\delta_{4 \rightarrow 3}(x_3)$
2	$\sum_{x_3} \psi(x_3) \psi(x_2, x_3) \delta_{4 \rightarrow 3}(x_3)$	$\delta_{3 \rightarrow 2}(x_2)$
3	$\sum_{x_2} \psi(x_2) \psi(x_2, x_1) \delta_{3 \rightarrow 2}(x_2)$	$\delta_{2 \rightarrow 1}(x_1)$

Note 1: Up to this point, VE & BP operations are IDENTICAL  
VE creates intermediate factors. } same  
BP calls them messages.

Note 2: At this point, VE stops & proceeds to compute  
 $P(x_i) \propto \psi(x_i) \delta_{2 \rightarrow 1}(x_i)$

But BP goes on in a "backward pass".

Note 3: After BP finishes backward pass of messages,  
we claim it has converged & no further  
message-passing is required. Why? Try writing  
out  $\delta_{4 \rightarrow 3}(x_3)$ . What do you find?

Key idea: leaves in trees only need to  
send messages ONCE. why?  
They don't multiply <sup>any</sup> incoming messages.  
So always done!

Note 4: After VE is done, it can only answer  
 $P(x_i | x_5 = x_5)$

After BP is done, it can answer

$$P(x_i | x_5 = x_5) \boxed{\forall i}! \quad \text{Cool!}$$
$$\propto \psi(x_i) \cdot \prod_{j \in \text{pa}(i)} \delta_{j \rightarrow i}(x_i)$$

At only 2x cost of VE, we get so much more.

2

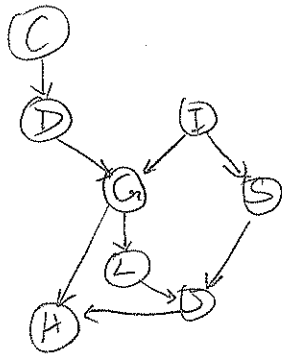
## ② Generalizing Connection of VE & BP beyond chains.

→ Key Idea: As VE is running on an arbitrary BN/MN, it creates something called a Clique Tree / Junction Tree.

This JT is in fact a valid cluster graph that BP can be run on.

Similar to chains, only need to run BP for 2 passes (Leaves → Root → leaves)

Example:



Elim Ordering: {C, D, I, H, G, S, L, J}

→ let's eliminate

→ Create clusters corresponding to factor multiplication

→ Add edges whenever an intermediate factor is "consumed" in another elimination

C:  $P(C)P(D|C)$

D:  $f_1(D)P(G|D,I)$

I:  $f_2(G,I)P(S|I)$

1: C D

2: D G I

3: G I S

$$\sum_C C(\cdot)C(\cdot) = f_1(D)$$

$$\sum_D C(\cdot)C(\cdot) = f_2(G,I)$$

$$\sum_I C(\cdot)C(\cdot) = f_3(G,S)$$

G:  $f_3(G,S)$   
 $f_4(G,I)$   
 $P(L|G)$

5: G J L S

H:  $P(H|G,S)$

4: G H J  
 $\sum_C C(\cdot) = f_4(G,S)$

$$\sum_G C(\cdot)C(\cdot) = f_5(J,L,S)$$

Note 1: Claim: This Junction Tree is a Cluster Graph  
∵ it satisfies (a) Fam Preserving Prop  
(b) RIP

Note 2: If we run BP with a leaves  $\rightarrow$  Root  $\rightarrow$  leaves schedule (as shown by directionality of edges) the first part leaves  $\rightarrow$  Root will be IDENTICAL to VE.

Note 3: <sup>\*</sup> Again, after backward pass BP will be converged.  
 $\Rightarrow$  Cluster graph calibrated  
 $\Rightarrow$  ready to answer queries  $P(x_i)$



### ③ Variational Inference:

$P(X)$ : Some complex distribution  
 $\rightarrow$  maybe BN, maybe MN  
 $\rightarrow$  maybe "reality"

$\rightarrow$  Want to answer  $P(x_i)$ . Too hard.

Key Idea:  $\rightarrow$  How about we approximate  $P(X)$  with a "simple" distribution  $Q(x)$  that is "close" to  $P(X)$ , but allows easy computation of  $Q(x_i)$ .

Specifically,  
Our goal

(3)

$$P(\vec{x}) = \frac{1}{Z} \prod_c \psi_c(x_c) \quad \text{MRF}$$

↑  
Notation means variables involved  
in  $c$ 'th clique

e.g. =  $\frac{1}{Z} \prod_i \psi_i(x_i) \prod_{\substack{(i,j) \\ \in E}} \psi_{ij}(x_i, x_j)$

[Pairwise MRF]

$P$  is too complex!

How about we approximate  $P$  with something <sup>VERY VERY</sup> simple

$Q$ : a fully factorized distribution

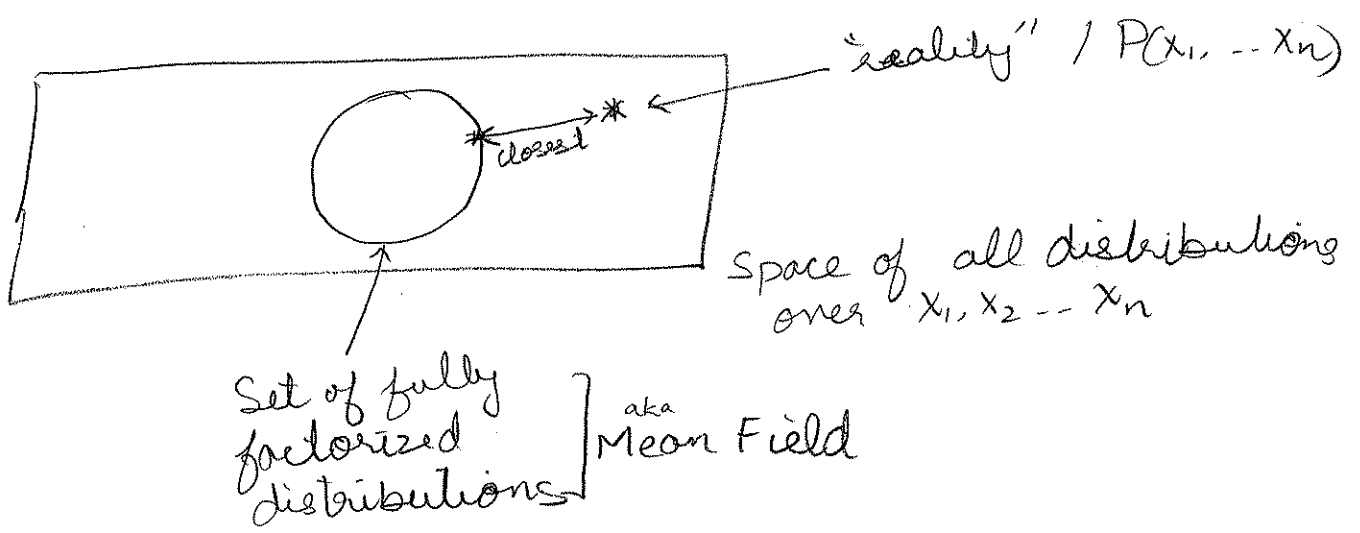
$$= \prod_i Q_i(x_i)$$

Once we find the  $Q$  closest to  $P$ , it is trivial to find marginals. Just read  $Q_i(x_i)$

want pairwise marginal?

No problem!  $Q_{ij}(x_i, x_j) = Q_i(x_i) Q_j(x_j)$

OK. So how do we find a good  $Q$ ?



Minimize a measure of "distance" between  $P$  &  $Q$

KL-divergence!

Two possibilities  $\rightarrow KL(P \parallel Q) = \sum_{\vec{x}} P(\vec{x}) \log \frac{P(\vec{x})}{Q(\vec{x})}$

$\because$  KL is not symmetric  $\rightarrow KL(Q \parallel P) = \sum_{\vec{x}} Q(\vec{x}) \log \frac{Q(\vec{x})}{P(\vec{x})}$

Note 1:  $\min_Q KL(P \parallel Q)$  is the "correct" direction.

Why?  $\because$  it weights  $\log \frac{P(\vec{x})}{Q(\vec{x})}$  with  $P(\vec{x})$  — the "true" prob of this event  $\vec{x}$ .

The other direction gets to ignore certain events by setting  $q(\vec{x}) \ll 1$  where  $q$  is a bad approximation of  $P$ .

See examples in slides.

4

④ Back to Graphical Models:

Variational Inference: the correct direction

$$P(\vec{x}) = \frac{1}{Z} \prod_c \psi_c(x_c)$$

$$q(\vec{x}) = \prod_i q_i(x_i)$$

Goal:  $\arg \min_q KL(P \parallel q) = \sum_{\vec{x}} P(\vec{x}) \log P(\vec{x}) - \sum_{\vec{x}} P(\vec{x}) \log q(\vec{x})$   
constant w.r.t.  $q$

$$= \arg \max_q \sum_{\vec{x}} P(\vec{x}) \log q(\vec{x})$$

$$= \sum_{\vec{x}} P(\vec{x}) \log \prod_i q_i(x_i)$$

$$= \sum_{\vec{x}} P(\vec{x}) \sum_i \log q_i(x_i)$$

$$= \sum_i \sum_{\vec{x}} P(\vec{x}) \log q_i(x_i)$$

$$= \sum_i \sum_{x_i} P_i(x_i) \log q_i(x_i)$$

Thus Optimization problem

$$\arg \max_{\{q_i\}} \sum_i \sum_{x_i} P_i(x_i) \log q_i(x_i)$$

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

$$q_i(x_i) \geq 0$$

Write Lagrangian  $L(\{q_i\}, \{\lambda_i\}) = \sum_i \sum_{x_i} P_i(x_i) \log q_i(x_i) - \sum_i \lambda_i \left[ \sum_{x_i} q_i(x_i) - 1 \right]$

Now set  $\frac{\partial L}{\partial q_i(x_i)} = 0$

$$\left\{ \frac{P_i(x_i)}{q_i(x_i)} - x_i = 0 \right.$$

[subtle point: assuming  $q_i(x_i) \neq 0$ ]

$$\Rightarrow q_i(x_i) = \frac{P_i(x_i)}{x_i}$$

$$\sum q_i(x_i) = 1 = \sum P_i(x_i)$$

$$\Rightarrow x_i^* = 1$$

Thus  $q_i^*(x_i) = P_i(x_i)$

Doesn't quite help!

Need to know marginals to compute  $Q$ .

But we were computing  $Q$  to get marginals!