$\Psi_{ij}(x_i, x_j)$



Deterministic Algorithms

$$P(\lbrace x_i\rbrace) = \frac{1}{Z} \prod_{i} \psi_i(x_i) \prod_{i,j} \psi_{ij}(x_i, x_j) \quad \text{MRF}$$

$$\psi_{ij}(x_i, x_j) \quad \text{Suppose, we want to estimate the marginal distributions } \prod_{i} P_i(x_i, x_i)$$

$$\psi_{ik}(x_i, x_k) \quad \psi_{jk}(x_j, x_k) \quad \text{or } \mathbf{x}^* = \arg \max_{\lbrace x_i\rbrace} P(\lbrace x_i\rbrace)$$

If the graph is a polytree (i.e. no closed loops), then we can use dynamic programming

 \rightarrow This cannot be applied if the graph has closed loops.

Here we describe a popular algorithm – **belief propagation** – that gives correct results on polytrees, and empirically good approximations most of time on graphs And we will show the relation to MCMC



Belief Propagation (BP)

proceeds by passing messages between the graph nodes

 $m_{ij}(x_j:t)$:message that node *i* passes to node *j* to affect state x_j

The messages gets updated as follows:

$$m_{ij}(x_j:t+1) = \sum_{x_i} \psi_{ij}(x_i, x_j) \psi_i(x_i) \prod_{k \neq j} m_{ki}(x_i:t) \quad \text{Sum-product rule}$$

Alternative: the max-product

$$m_{ij}(x_j:t+1) = \max_{x_i} \left\{ \psi_{ij}(x_i, x_j) \psi_i(x_i) \prod_{k \neq j} m_{ki}(x_i:t) \right\}$$

If the algorithm converges (it may not), then we compute approximations to the marginals: $b_i(x_i) \propto \psi_i(x_i) \prod_k m_{ki}(x_i)$ $b_{ij}(x_i, x_j) \propto \psi_i(x_i) \psi_{ij}(x_i, x_j) \psi_j(x_j) \prod_{k \neq j} m_{ki}(x_i) \prod_{l \neq i} m_{lj}(x_j)$ Let

Lecture BP-02



Belief Propagation

BP (sum-product) was first proposed by Judea Pearl (C.S. UCLA) for performing inference on polytrees

The max-product algorithm was proposed earlier by Gallagher

The application was developed in the 1990's for decoding problems

- goal was to achieve Shannon's bound on information transmission

Experimentally, it was shown that BP usually converges to reasonable approximations

- Full understanding of when & why it converges is an open problem

On polytree, it is similar to dynamic programming – so in a sense, it is the way to extend DP to graphs with closed loops



Example of BP (sum-product)

$$\mathbf{x} = (x_1, x_2, x_3, x_4)$$

$$P(\mathbf{x}) = \frac{1}{Z} \psi_{12}(x_1, x_2) \psi_{23}(x_2, x_3) \psi_{34}(x_3, x_4)$$

Messages:

 $m_{12}(x_2)$, from node 1 to node 2 $m_{21}(x_1)$, from node 2 to node 1 $m_{23}(x_3)$, from node 2 to node 3 $m_{32}(x_2)$, from node 3 to node 2 $m_{34}(x_4)$, from node 3 to node 4 $m_{43}(x_3)$, from node 4 to node 3

Update rule (from message passing equation) $m_{12}(x_2:t+1) = \sum_{x_1} \psi_{12}(x_1, x_2)$ → boundary $m_{21}(x_1:t+1) = \sum_{x_2} \psi_{12}(x_1, x_2) m_{32}(x_2:t)$ $m_{23}(x_3:t+1) = \sum_{x_2} \psi_{23}(x_2, x_3) m_{12}(x_2:t)$ $m_{32}(x_2:t+1) = \sum_{x_3} \psi_{23}(x_2, x_3) m_{43}(x_3:t)$ $m_{34}(x_4:t+1) = \sum_{x_3} \psi_{34}(x_3, x_4) m_{23}(x_3:t)$ $m_{43}(x_3:t+1) = \sum_{x_4} \psi_{34}(x_3, x_4)$ → boundary Lecture BP-04



Many ways to run BP

(1) Update all messages in parallel

Note: BP can be parallelized but Dynamic Programming cannot)

(2) Start at boundaries



Calculate $m_{12}(x_2)$, then $m_{23}(x_3)$, then $m_{34}(x_4)$ Forward pass (like Dynamic Programming)

> $m_{43}(x_3)$, then $m_{32}(x_2)$, then $m_{21}(x_1)$ \rightarrow Backward pass (like backward pass of DP)

Then read off estimates of unary marginals $b_{1}(x_{1}) = \frac{1}{Z_{1}} m_{21}(x_{1}), \quad Z_{1} = \sum_{x_{1}} m_{21}(x_{1}) :\text{normalization}$ $b_{2}(x_{2}) = \frac{1}{Z_{2}} m_{12}(x_{2}) m_{32}(x_{2}), \quad Z_{2} :\text{normalization}$ $b_{3}(x_{3}) = \frac{1}{Z_{3}} m_{23}(x_{3}) m_{43}(x_{3}), \quad Z_{3} :\text{normalization}$ $b_{4}(x_{4}) = \frac{1}{Z_{4}} m_{34}(x_{4}), \quad Z_{4} :\text{normalization}$ $b_{12}(x_{1}, x_{2}) = \frac{1}{Z_{12}} \psi_{12}(x_{1}, x_{2}) m_{32}(x_{2}), \quad \text{and so on}$ Lecture BP-05



Many ways to run BP

Alternatively, initialize the m's to take an initial value – e.g. $m_{ij}(x_j)=1$ for all i, jand update the messages in any order will still converge for graph with no closed loops Then estimates (beliefs) will be the true marginals

e.g.
$$b_1(x_1) = \sum_{x_2, x_3, x_4} P(x_1, x_2, x_3, x_4)$$

 $b_{12}(x_1, x_2) = \sum_{x_3, x_4} P(x_1, x_2, x_3, x_4)$

But, BP will often converge to a good approximation to the marginals for graphs which do not have closed loops This was discovered in the late 1980's



Advantages of BP over DP

(1) BP will converges (approximates) for many graphs with closed loops

(2) BP is parallelizable (nice if you have a parallel computer or GPU)

An alternative way to consider BP

Make local approximations to the local distribution (BP w/o messages)

$$B(x_{i}, \mathbf{x}_{N(i)}) = \frac{1}{Z} b_{i}(x_{i}) \prod_{j \in N(i)} \frac{b_{ij}(x_{i}, x_{j})}{b_{i}(x_{i})}$$

If the $b_{ij}(x_i, x_j) \& b_i(x_i)$ are the true marginal distribution, then

$$\frac{b_{ij}(x_i, x_j)}{b_i(x_i)} = b(x_j \mid x_i)$$

Lecture BP-07

An alternative way to consider BP



Lecture BP-08



How does this relate to MCMC?

Recall that Gibbs Sampler $K_{\gamma}(\mathbf{x} | \mathbf{x}') = P(x_{\gamma} | x'_{N(\gamma)})S_{\mathbf{x}/\gamma,\mathbf{x}'/\gamma}$

Substituting the Gibbs sampler into the Chapman-Kolmogorov equation $\mu_{t+1}(\mathbf{x}_{\gamma}) = \sum_{\mathbf{x}_{N(\gamma)}} P(\mathbf{x}_{r} | \mathbf{x}'_{N(\gamma)}) \mu_{t}(\mathbf{x}_{N(\gamma)})$ Replace $\mu_{t+1}(\mathbf{x}_{\gamma})$ by $\sum_{x_{i}} B(x_{i}, x_{N(\gamma)}), \quad \mathbf{x}_{r} = y_{i}$ or $\sum_{x_{i}, x_{j}} B(x_{i}, x_{j}, x_{N(i,j)}), \quad \mathbf{x}_{\gamma} = (x_{i}, x_{j})$ Lecture BP-09



Bethe Free Energy

It can be shown that the fixed point of BP correspond to extreme of the Bethe Free Energy

$$F[b] = \sum_{ij} \sum_{x_i, x_j} b_{ij}(x_i, x_j) \ln \frac{b_{ij}(x_i, x_j)}{\psi_i(x_i)\psi_j(x_j)\psi_{ij}(x_i, x_j)} - \sum_i (n_i - 1) \sum_{x_i} b_i(x_i) \ln \frac{b_i(x_i)}{\psi_i(x_i)}$$

This leads an alternative class of algorithm which seek to directly minimize F[b]These algorithms are more complex than BP, more time consuming, and do not always give better results

Note M. Wainwright defines a class of convex free energies similar to Bethe **Note** Junction trees allows DP to be applied to same graphs with closed

loops (see Lauritzen and Spiegelhalter). Lecture BP-10



A range of alternative algorithms

The original is meanfield (MFT) Kulback-Leibler: Define $B(\mathbf{x}) = \prod_i b_i(x_i)$ $KL(B) = \sum_{\mathbf{x}} B(\mathbf{x}) \ln \frac{B(\mathbf{x})}{P(\mathbf{x})}$ Seek to find the $B(\mathbf{x})$ that minimizes KL(B)Equivalent to $\sum_{i,j} \sum_{x_i,x_j} b_i(x_i) b_j(x_j) \ln \frac{b_i(x_i)b_j(x_j)}{\psi_i(x_i)\psi_{ij}(x_i,x_j)\psi_j(x_j)} - \sum_i (n_i - 1) \sum_{x_i} b_i(x_i) \ln \frac{b_i(x_i)}{\psi_i(x_i)}$ Compare to Bethe Free Energy: $b_{ij}(x_i, x_j) \rightarrow b_i(x_i) b_j(x_j)$ Minimizing KL(B) is not straightforward, but it is straightforward to find a local minima

• These approaches are significantly faster than MCMC, but MCMC works when these do not.