

Approximate inference using loopy belief propagation

Yair Weiss

Hebrew University of Jerusalem

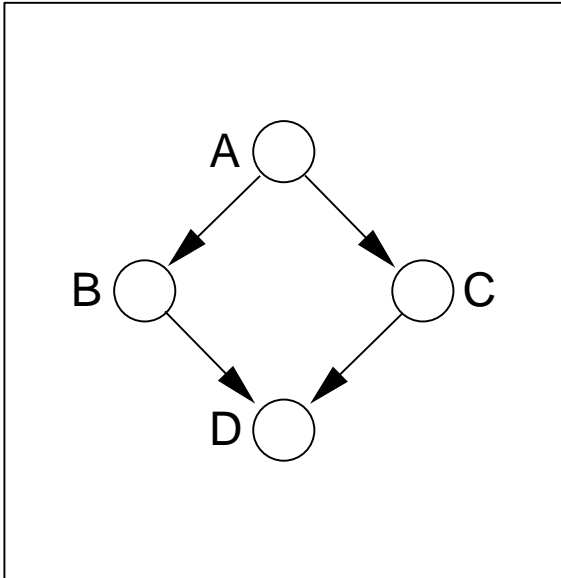
www.cs.huji.ac.il/~yweiss

Outline

- Exact inference using belief propagation.
- The need for approximation.
- Empirical success of loopy BP (and some failures).
- Theoretical results on loopy BP.
- Beyond ordinary BP.

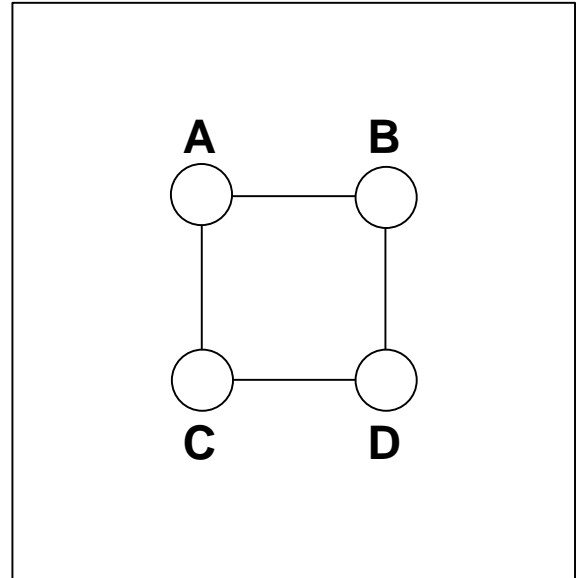
Graphical Models

Bayes nets



$$\prod_x P(x|\text{par}(x))$$

Markov Nets

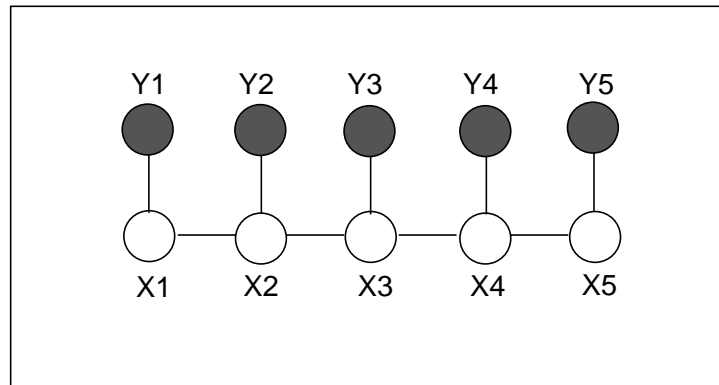


$$\prod_c \Psi(x_c)$$

Nodes: random variables.

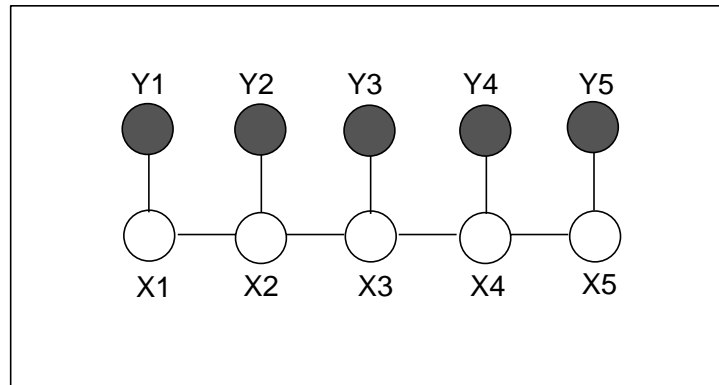
Edges: probabilistic constraints between variables.

Inference in Graphical Models



1. Marginals. Find $P(X_i = x_i|Y)$.
2. MAP assign. Find $\{x_i^*\}$ such that $P(\{X_i = x_i^*\}|Y)$ is max.
3. MM assign. Find $\{x_i^*\}$ such that $P(X_i = x_i^*|Y)$ is max for each i .
4. Likelihood. Calculate $P(Y)$.

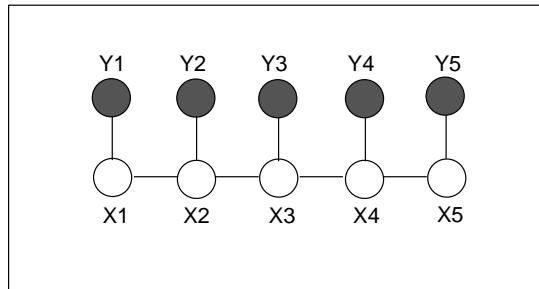
The forward-backward algorithm for HMMs



$$\begin{aligned}P(x_t|y) &\propto \alpha_t(x_t)\beta_t(x_t)P(y_t|x_t) \\ \alpha_t(x_t) &= P(x_t, y_{1,t-1}) \\ \beta_t(x_t) &= P(y_{t+1,T}|x_t)\end{aligned}$$

$$\begin{aligned}\beta_t(x_t) &= P(y_{t+1,T}|x_t) \\ &= \sum_{x_{t+1}} P(x_{t+1}, y_{t+1,T}|x_t) \\ &= \sum_{x_{t+1}} P(x_{t+1}|x_t)P(y_{t+1}|x_{t+1})\beta_{t+1}(x_{t+1}) \\ \alpha_t(x_t) &= \sum_{x_{t-1}} P(x_t|x_{t-1})P(y_{t-1}|x_{t-1})\alpha_{t-1}(x_{t-1})\end{aligned}$$

Belief Propagation for undirected trees



Input: Graph, $\Psi_{ij}(x_i, x_j)$, $\Psi_{ii}(x_i, y_i)$.

Define:

$m_{ij}(x_j)$ message that X_i sends to X_j .

$b_i(x_i)$ "belief" at node X_i .

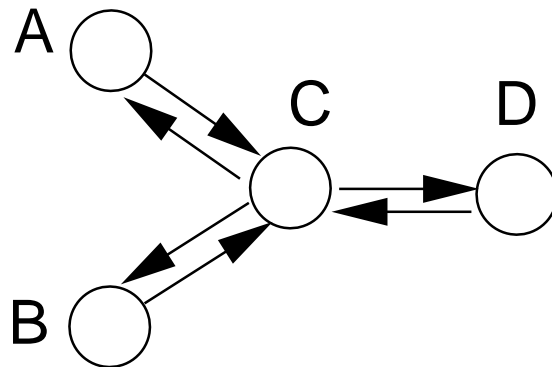
Iterate:

$$m_{ij}(x_j) \leftarrow \alpha \sum_{x_i} \Psi_{ij}(x_i, x_j) \Psi_i(x_i) \prod_{X_k \in N(X_i) \setminus X_j} m_{ki}(x_i)$$

finally:

$$b_i(x_i) \leftarrow \alpha \Psi_i(x_i) \prod_{X_j \in N(X_i)} m_{ji}(x_i)$$

BP Dynamics:

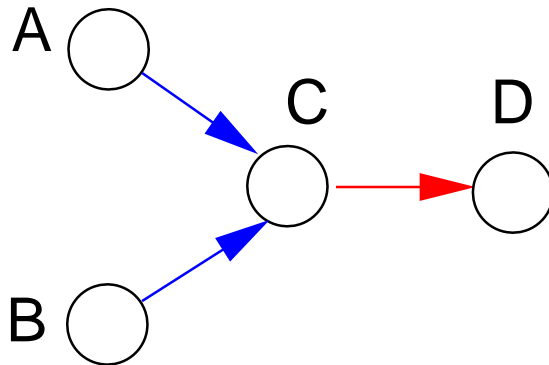


$$m_{ij}(x_j) \leftarrow \alpha \sum_{x_i} \psi_{ij}(x_i, x_j) \psi_i(x_i) \prod_{X_k \in N(X_i) \setminus X_j} m_{ki}(x_i)$$

A parallel message-passing algorithm.

- Every node sends a probability density to its neighbors.
- message to neighbor depends on messages received from all the *other neighbors*

BP Dynamics

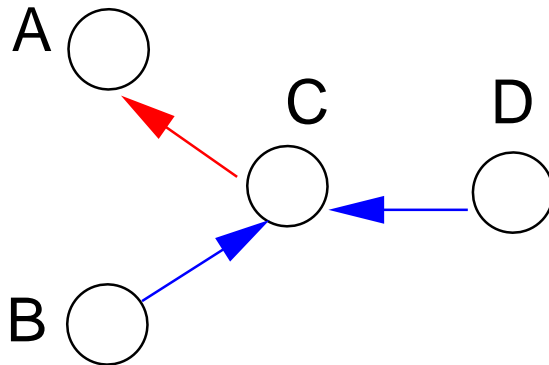


$$m_{ij}(x_j) \leftarrow \alpha \sum_{x_i} \Psi_{ij}(x_i, x_j) \Psi_i(x_i) \prod_{X_k \in N(X_i) \setminus X_j} m_{ki}(x_i)$$

A parallel message-passing algorithm.

- Every node sends a probability density to its neighbors.
- message to neighbor depends on messages received from all the *other neighbors*

BP Dynamics:

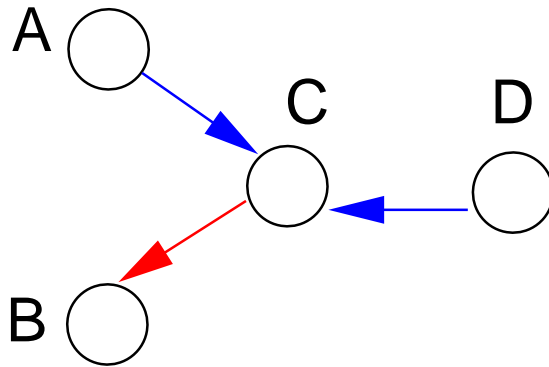


$$m_{ij}(x_j) \leftarrow \alpha \sum_{x_i} \psi_{ij}(x_i, x_j) \psi_i(x_i) \prod_{X_k \in N(X_i) \setminus X_j} m_{ki}(x_i)$$

A parallel message-passing algorithm.

- Every node sends a probability density to its neighbors.
- message to neighbor depends on messages received from all the *other neighbors*

BP Dynamics

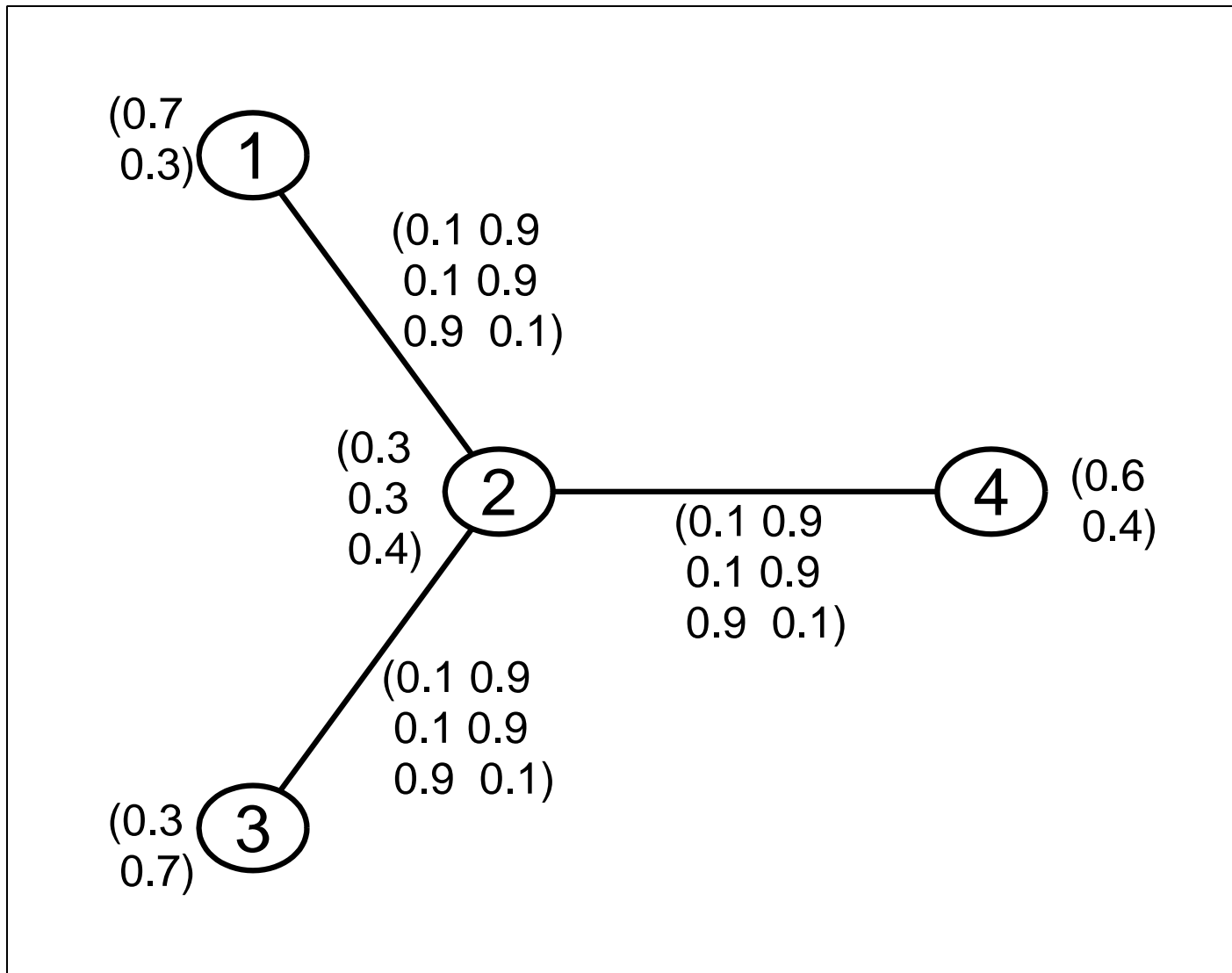


$$m_{ij}(x_j) \leftarrow \alpha \sum_{x_i} \psi_{ij}(x_i, x_j) \psi_i(x_i) \prod_{X_k \in N(X_i) \setminus X_j} m_{ki}(x_i)$$

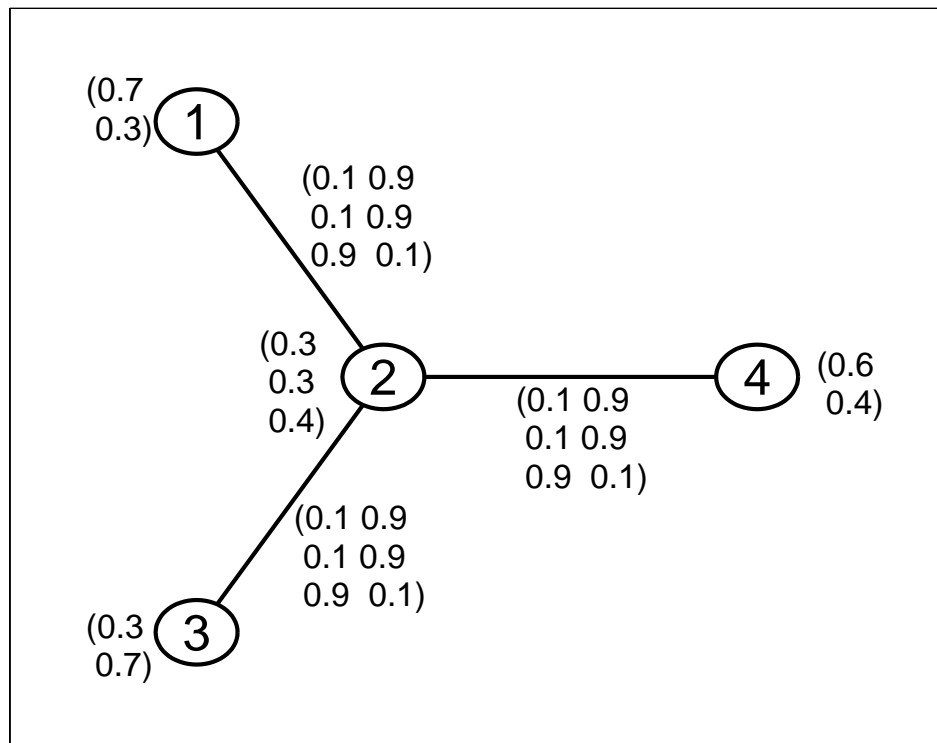
A parallel message-passing algorithm.

- Every node sends a probability density to its neighbors.
- message to neighbor depends on messages received from all the *other neighbors*

Numerical Example:



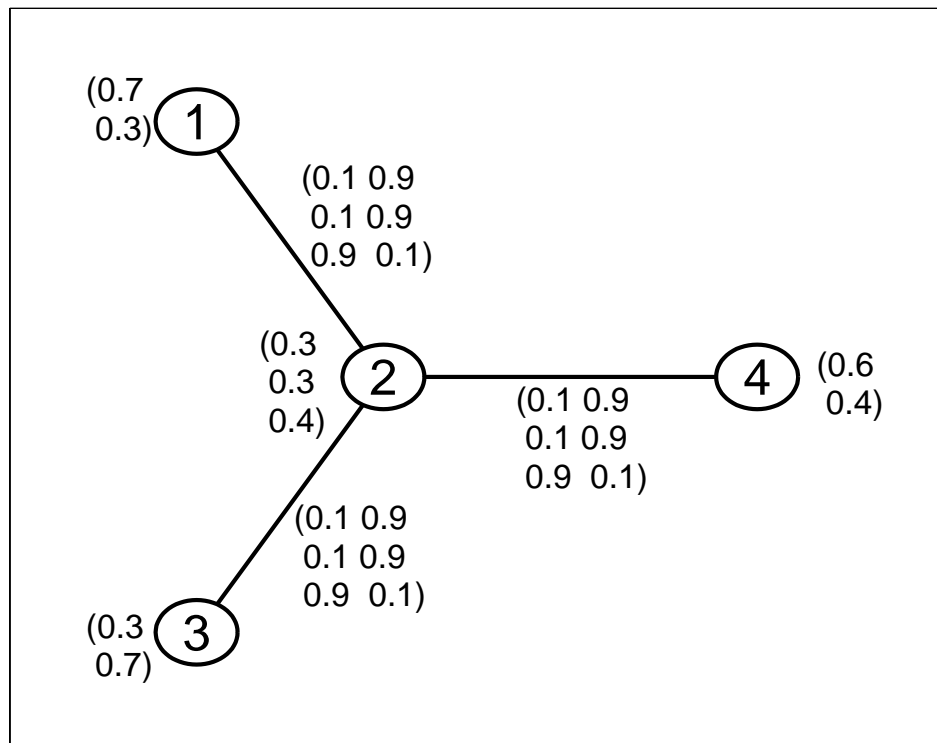
Numerical Example (cont):



$$m_{42} = \alpha \begin{pmatrix} 0.1 & 0.9 \\ 0.1 & 0.9 \\ 0.9 & 0.1 \end{pmatrix} \begin{pmatrix} 0.6 \\ 0.4 \end{pmatrix} = \begin{pmatrix} 0.2958 \\ 0.2958 \\ 0.4085 \end{pmatrix}$$

$$m_{32} = \alpha \begin{pmatrix} 0.1 & 0.9 \\ 0.1 & 0.9 \\ 0.9 & 0.1 \end{pmatrix} \begin{pmatrix} 0.3 \\ 0.7 \end{pmatrix} = \begin{pmatrix} 0.3976 \\ 0.3976 \\ 0.2048 \end{pmatrix}$$

Numerical Example (cont):

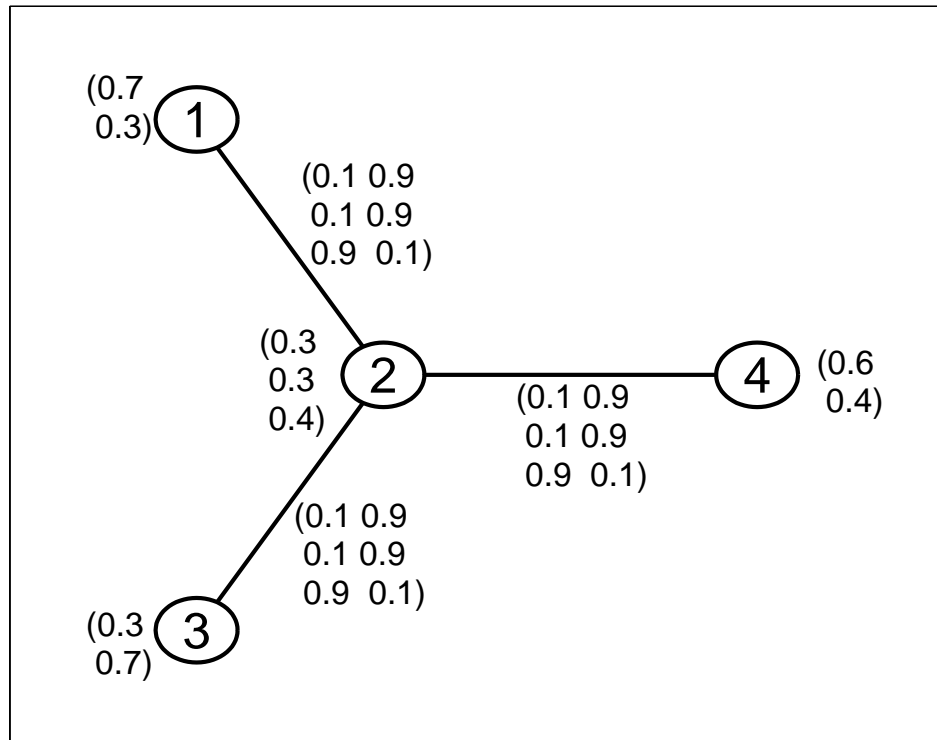


$$m_{21} =$$

$$\begin{pmatrix} 0.1 & 0.1 & 0.9 \\ 0.9 & 0.9 & 0.1 \end{pmatrix} \begin{pmatrix} 0.29 \\ 0.29 \\ 0.40 \end{pmatrix} * \begin{pmatrix} 0.39 \\ 0.39 \\ 0.20 \end{pmatrix} * \begin{pmatrix} 0.3 \\ 0.3 \\ 0.4 \end{pmatrix}$$

$$= \begin{pmatrix} 0.2869 \\ 0.7131 \end{pmatrix}$$

Numerical Example (cont):



$$b_1 = \alpha \begin{pmatrix} 0.2869 \\ 0.7131 \end{pmatrix} * \begin{pmatrix} 0.7 \\ 0.3 \end{pmatrix} = \begin{pmatrix} 0.484 \\ 0.516 \end{pmatrix}$$

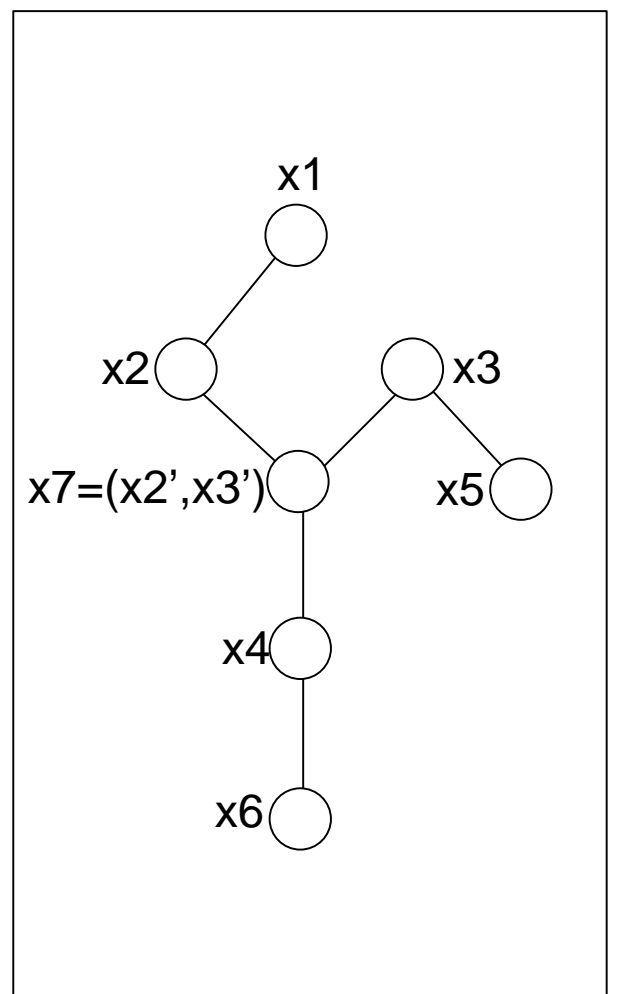
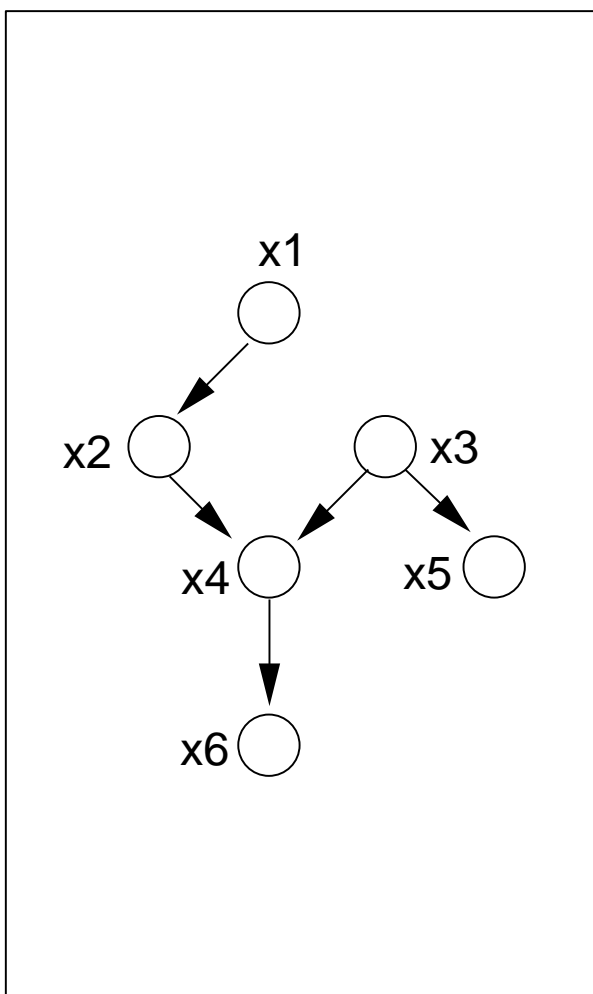
$$b_2 = \alpha \begin{pmatrix} 0.29 \\ 0.29 \\ 0.40 \end{pmatrix} * \begin{pmatrix} 0.39 \\ 0.39 \\ 0.20 \end{pmatrix} * \begin{pmatrix} 0.253 \\ 0.253 \\ 0.492 \end{pmatrix} * \begin{pmatrix} 0.3 \\ 0.3 \\ 0.4 \end{pmatrix}$$

$$= \begin{pmatrix} 0.253 \\ 0.253 \\ 0.492 \end{pmatrix}$$

Belief propagation in directed polytrees

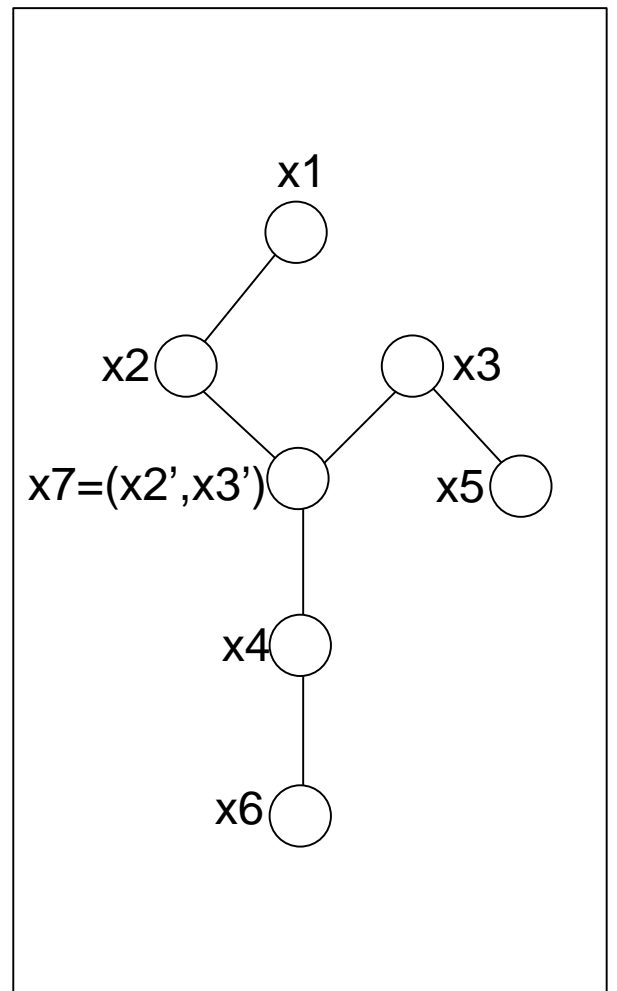
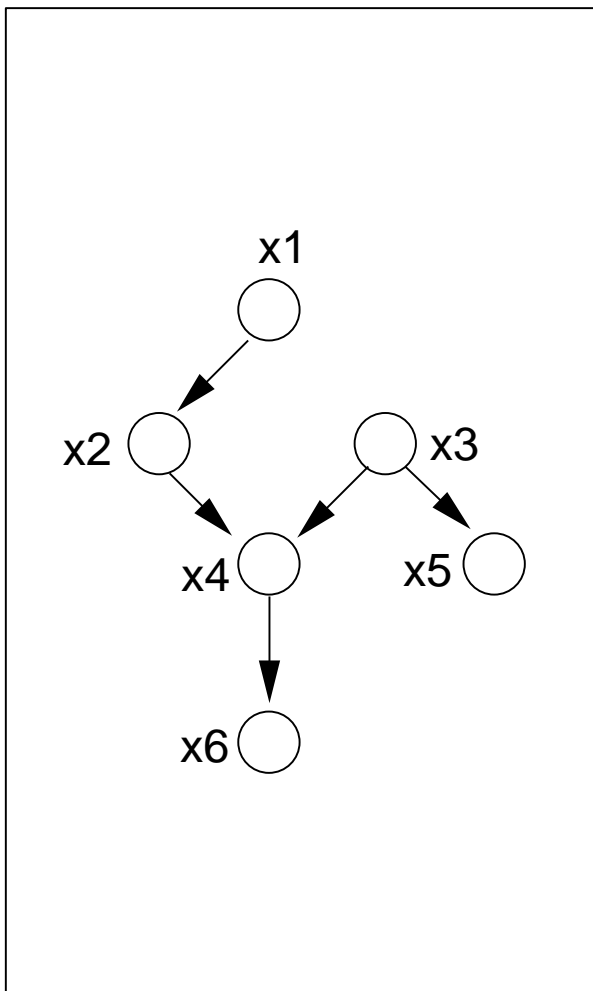
Many equivalent formulations (Including: Kim and Pearl 83, Jensen 94, Shafer and Shenoy 90, Aji and McEliece 97, Frey et al 97).

Equivalent to converting to undirected graph and running undirected BP.



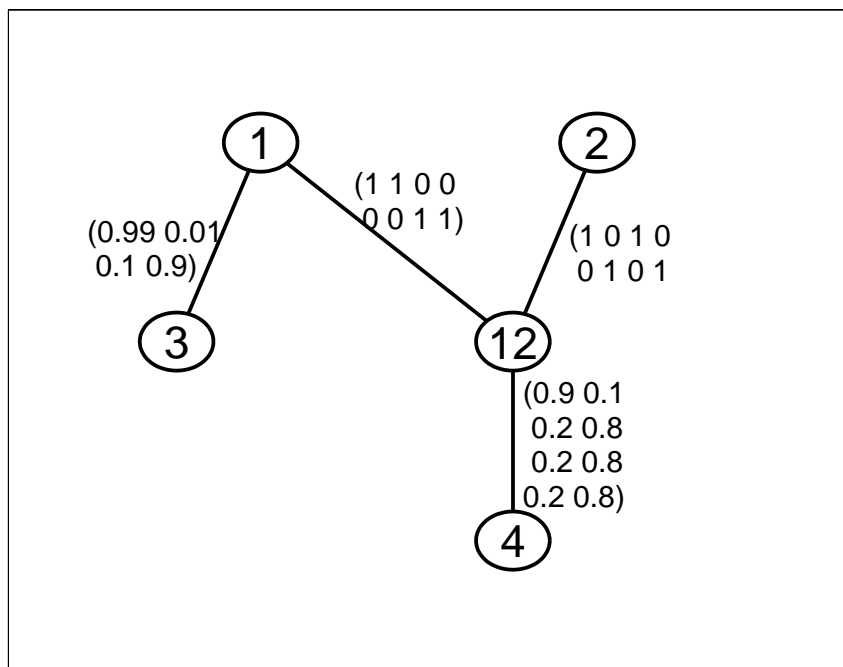
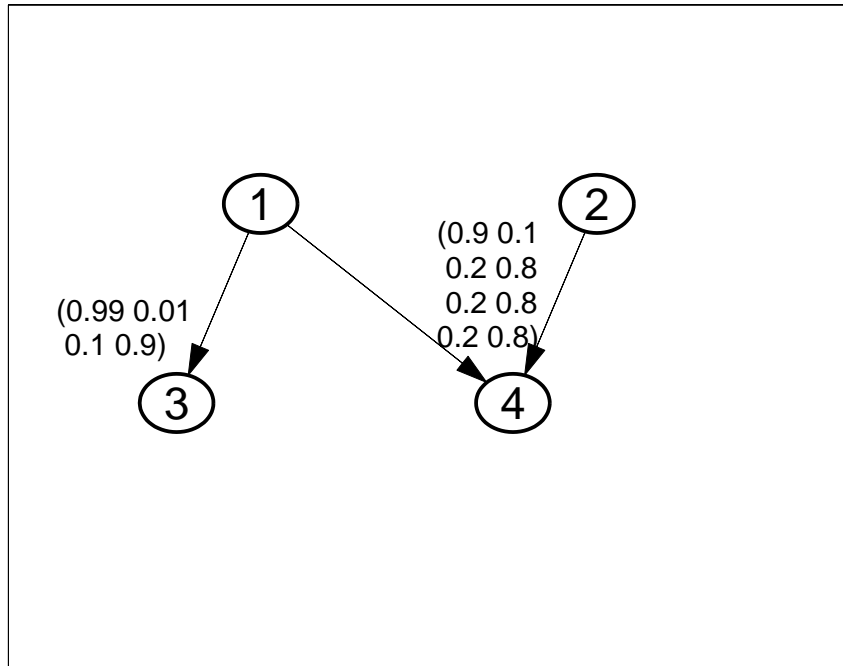
Belief propagation in directed polytrees

$$\begin{aligned}\psi_{46}(x_4, x_6) &= P(x_6|x_4) \\ \psi_{74}(x'_2, x'_3, x_4) &= P(x_4|x'_2, x'_3) \\ \psi_{27}(x_2, x'_2, x'_3) &= \delta(x_2 - x'_2)\end{aligned}$$

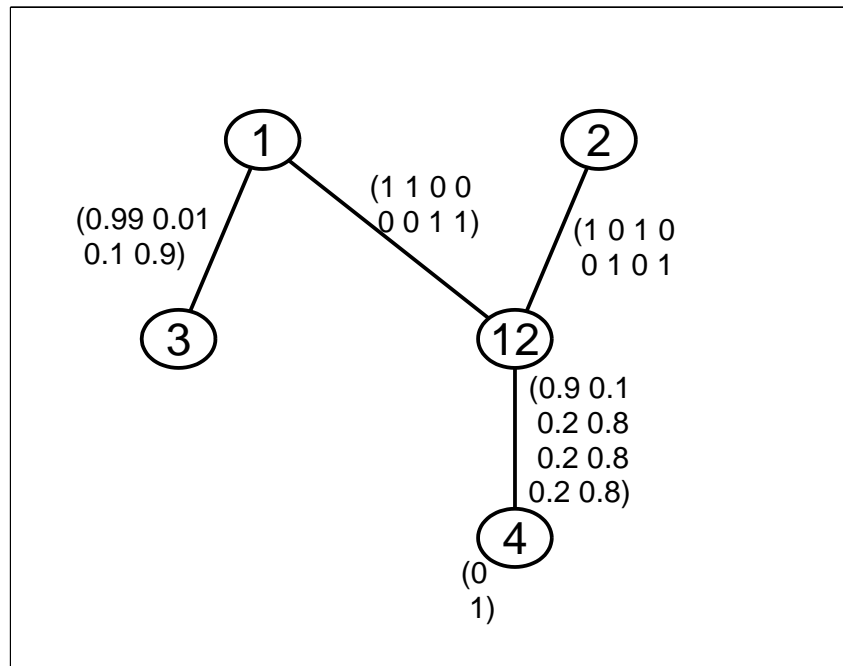


Example:

1:Earthquake,2:Burglary,3:Radio,4:Alarm.

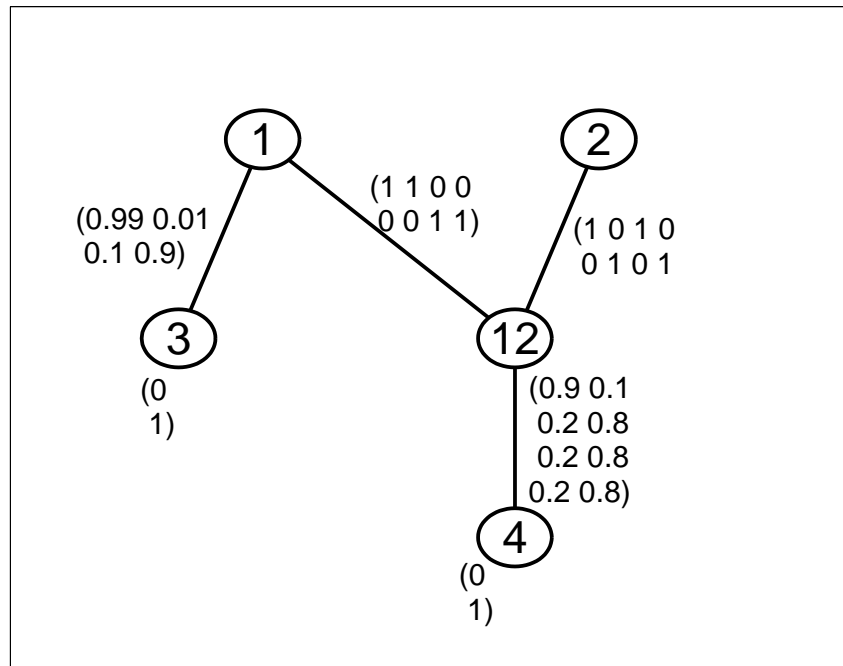


Numerical example (cont)



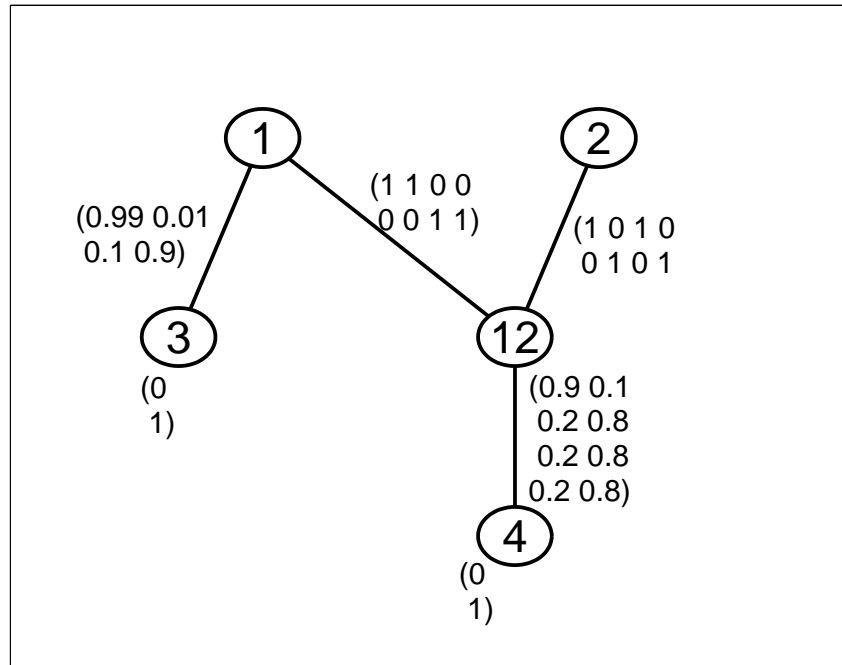
$$\begin{aligned}
 m_{12,2} &= \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} 0.9 & 0.1 \\ 0.2 & 0.8 \\ 0.2 & 0.8 \\ 0.2 & 0.8 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \end{pmatrix} \\
 &= \begin{pmatrix} 0.36 \\ 0.64 \end{pmatrix}
 \end{aligned}$$

Explaining away



$$\begin{aligned}
 m_{1,12} &= \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 0.99 & 0.01 \\ 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \end{pmatrix} \\
 &= \begin{pmatrix} 0.01 \\ 0.01 \\ 0.49 \\ 0.49 \end{pmatrix} \\
 \Rightarrow \\
 m_{12,2} &= \begin{pmatrix} 0.49 \\ 0.51 \end{pmatrix}
 \end{aligned}$$

Q: How is this related to Pearl's algorithm?



A: it's exactly the same!

$$m_{12,2}(x_2) = \lambda(x_2).$$

$$m_{12,4}(x_4) = \pi(x_4)$$

Two variants of belief propagation

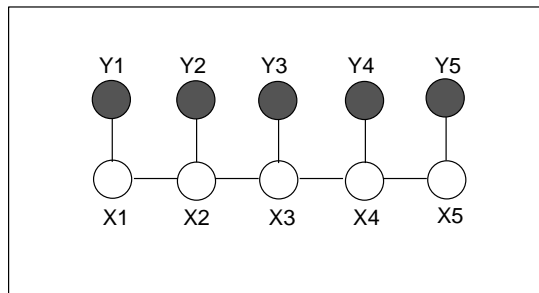
Sum-Product algorithm (a.k.a belief update)

$$m_{ij}(x_j) \leftarrow \alpha \sum_{x_i} \Psi_{ij}(x_i, x_j) \Psi_i(x_i) \prod_{X_k \in N(X_i) \setminus X_j} m_{ki}(x_i)$$

Max-Product algorithm (a.k.a belief revision, Viterbi algorithm, Dawid's algorithm)

$$m_{ij}(x_j) \leftarrow \alpha \max_{x_i} \Psi_{ij}(x_i, x_j) \Psi_i(x_i) \prod_{X_k \in N(X_i) \setminus X_j} m_{ki}(x_i)$$

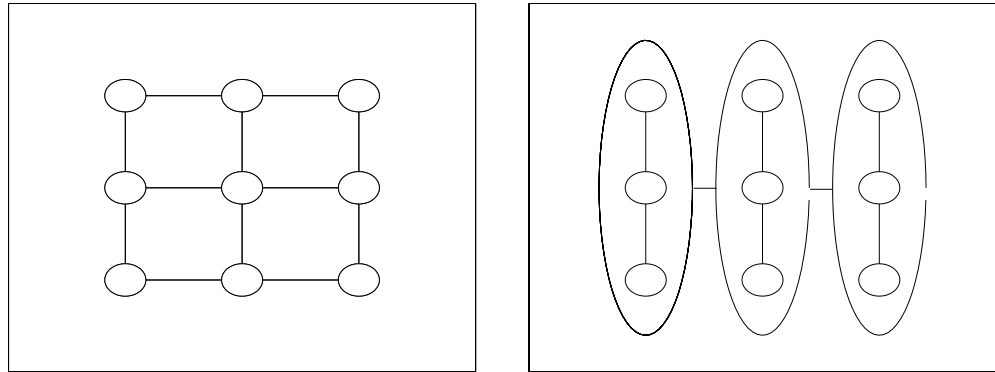
Belief Propagation in trees



When the graph is a tree (Pearl 86):

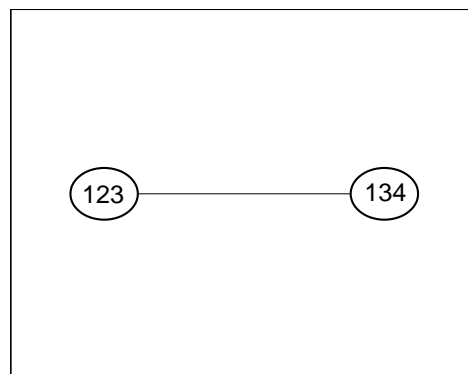
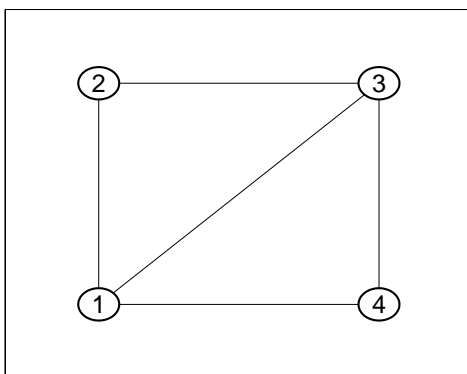
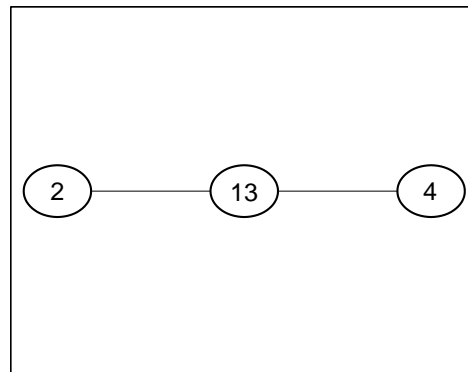
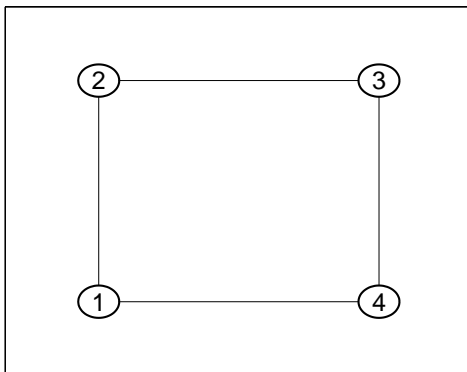
- both algorithms converge in finite time.
- at convergence: $b_i(x_i) = P(X_i = x_i | Y)$
- max-product assignment
($x_i^* = \max_{x_i} b_i(x_i)$ after convergence)
gives MAP assignment.
- sum-product assignment
($x_i^* = \max_{x_i} b_i(x_i)$ after convergence)
gives MM assignment.

Exact inference in arbitrary graphical model



- Convert graph to an equivalent tree.
- run BP on the tree.

Converting arbitrary graph to a tree



- clustering (e.g. Pearl 86).
- Junction tree (e.g. Lauritzen 96)
- Bucket tree (e.g. Dechter 97)

Optimal conversion NP hard.

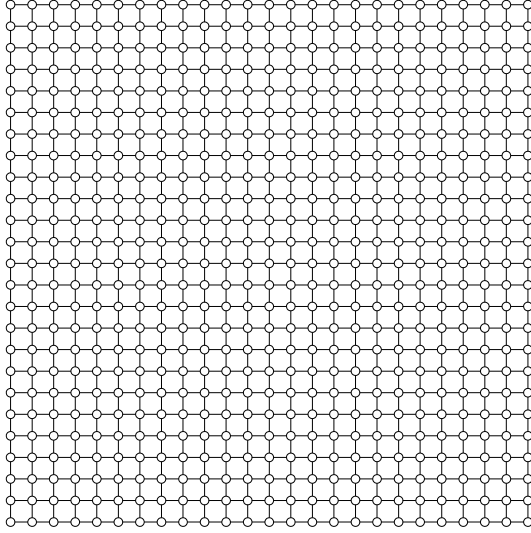
The need for approximate inference

We can do exact inference in any graphical model by converting to a tree and running BP, but:

- Complexity of BP is exponential in the size of the nodes *.
- In most real world applications, node size in equivalent tree is huge.

* some exceptions: noisy or CPTs, Gaussian CPTs.

Example: image analysis

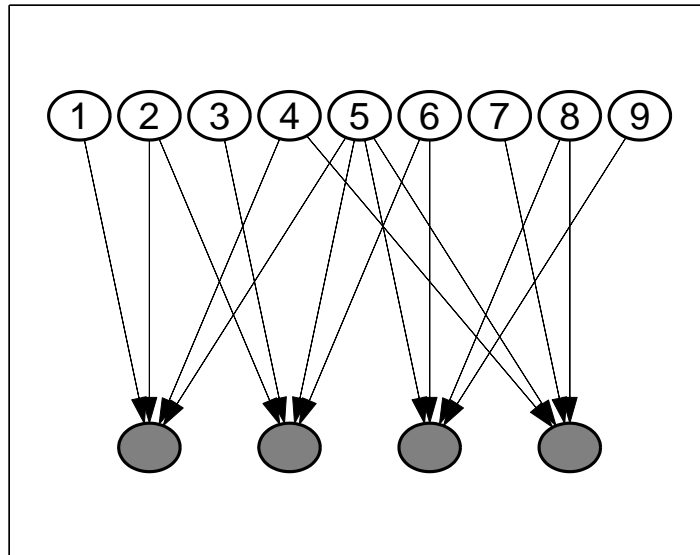


We want to classify pixels in an aerial photo as vegetation or not.

At each pixel we have color measurements and a prior that nearby pixels tend to have the same label.

Goal: calculate probability that a pixel is vegetation given measurements.

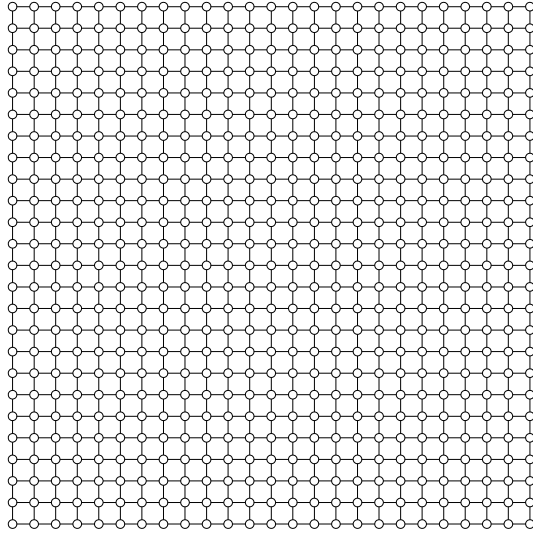
Example: LDPC Error correcting codes



Low power wireless communication \Rightarrow noise. We use codewords in which subsets of the bits are constrained to have zero parity.

Goal: calculate probability that a bit is a one.

The computational bottleneck



Assume the image size is 256x256 then *exactly* calculating the probability that a pixel is vegetation takes order 2^{256} operations.

Similar numbers for exact inference in LDPCs.

Approximate inference in graphical models

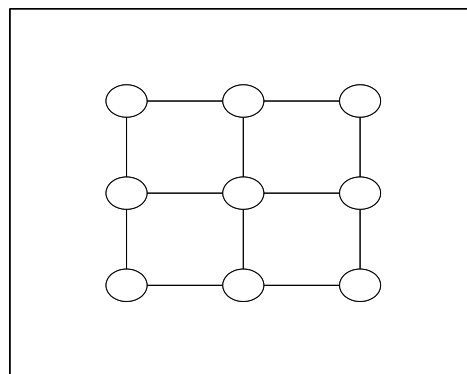
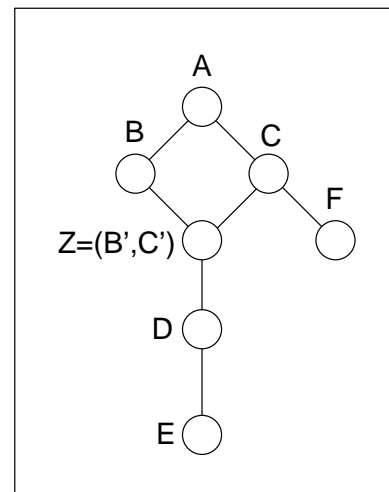
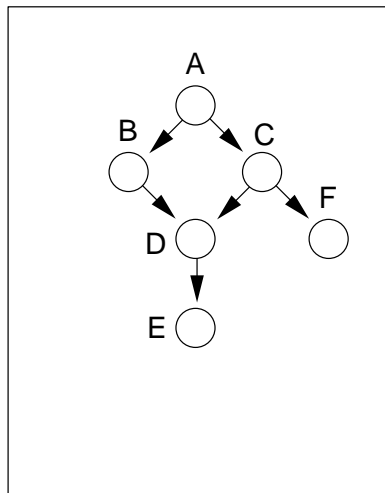
- Monte-Carlo methods (e.g. Geman and Geman 84)
- Variational (“mean-field”) approaches (e.g. Geiger and Girosi 91, Jordan et al. 98)
- “loopy” belief propagation.

Loopy Belief propagation in directed graphs

Convert graph to pairwise undirected graph.

Iterate, loopy BP equation:

$$m_{ij}(x_j) \leftarrow \alpha \sum_{x_i} \Psi_{ij}(x_i, x_j) \Psi_i(x_i) \prod_{X_k \in N(X_i) \setminus X_j} m_{ki}(x_i)$$



Loopy belief propagation

“When loops are present, the network is no longer singly connected and local propagation schemes will invariably run into trouble ... If we ignore the existence of loops and permit the nodes to continue communicating with each other as if the network were singly connected, messages may circulate indefinitely around the loops and the process may not converge to a stable equilibrium ... (even if it does) this asymptotic equilibrium is not coherent, in the sense that it does not represent the posterior probabilities of all nodes of the network” (Pearl 1988, p. 195)

Empirical successes of loopy propagation

- Gallager 63. Berrou et al. 93. Error correcting codes.
- Freeman and Pasztor 98, Saund 98. Image understanding.
- Shazeer et al. 99. Automatic crossword puzzle solution.
- Murphy et al. 99. Medical diagnosis.
- Frey 99. Factor analysis.
- Fridman and Mumford 99. Ising model.
- Frey 2001. Phase Unwrapping.

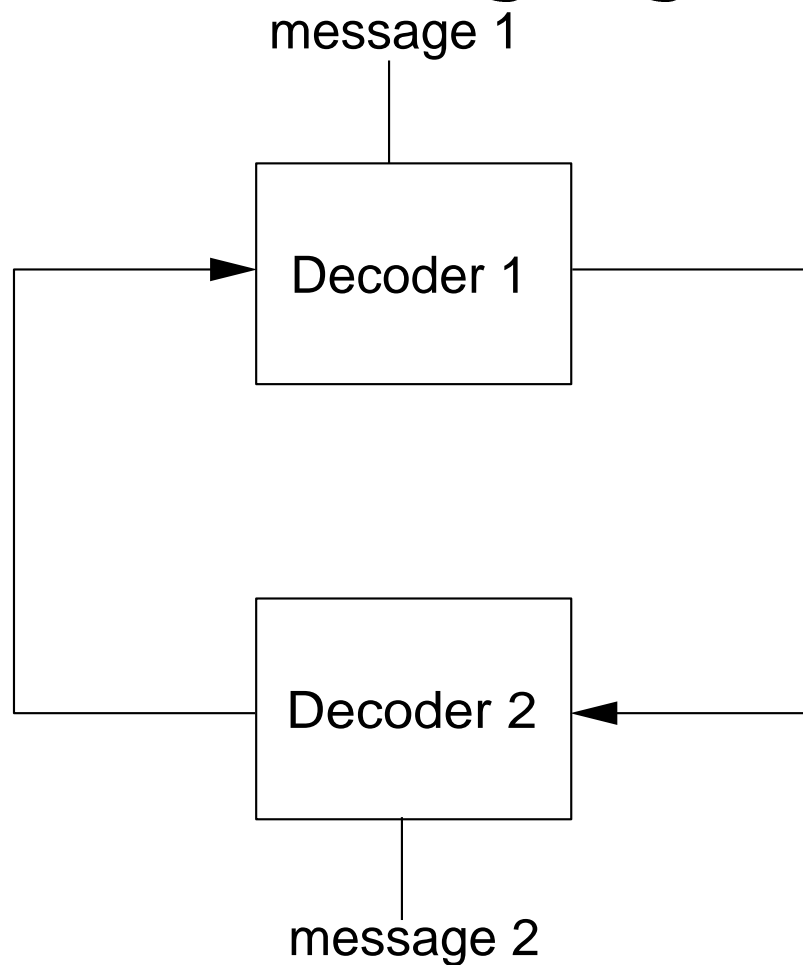
Error-correcting codes

X_i are bits. Original bits plus redundant “check” bits. Y_i are these bits sent over a noisy channel.

We want to design a code that:

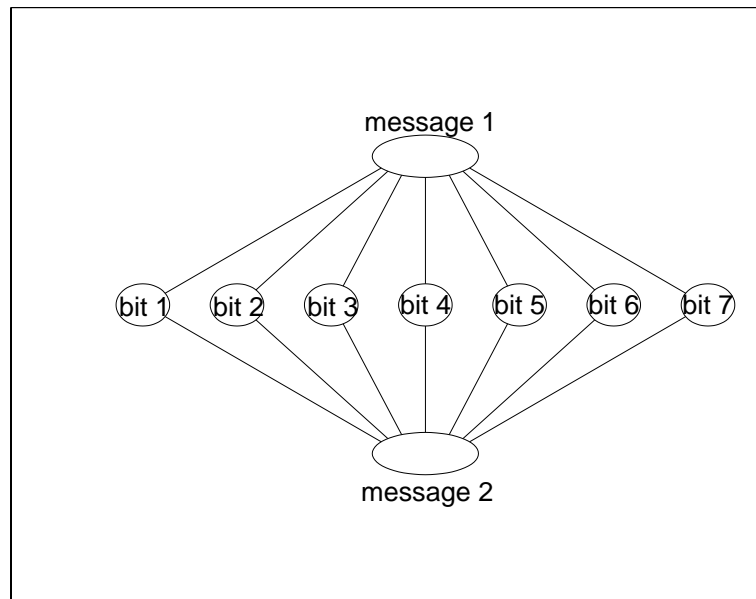
- can be decoded in reasonable time.
- the probability of error is as low as possible (Shannon’s limit).

Turbo decoding algorithm



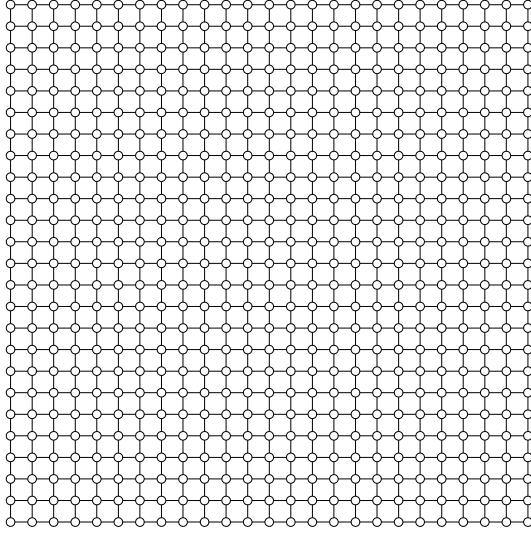
“A genuine and perhaps historic breakthrough”
... “the most exciting and potentially important
development in coding theory in many years”
(McElice et al 95)

Turbo decoding is belief propagation



(McEliece et al 97, Wiberg 96). Belief propagation is not supposed to work in this case.

Example: image analysis

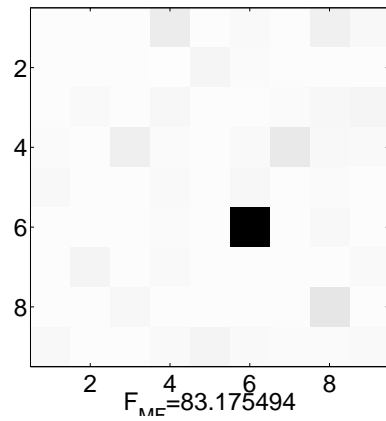


We want to classify pixels in an aerial photo as vegetation or not.

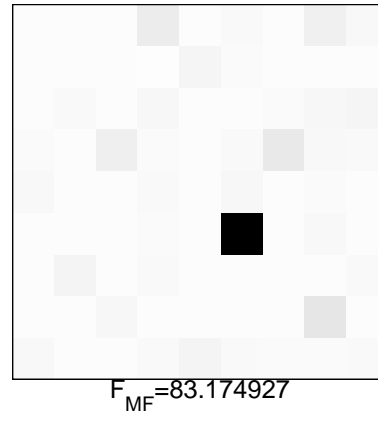
At each pixel we have color measurements and a prior that nearby pixels tend to have the same label.

Goal: calculate probability that a pixel is vegetation given measurements.

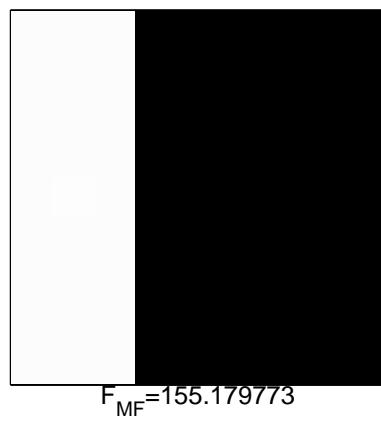
approximate inference results



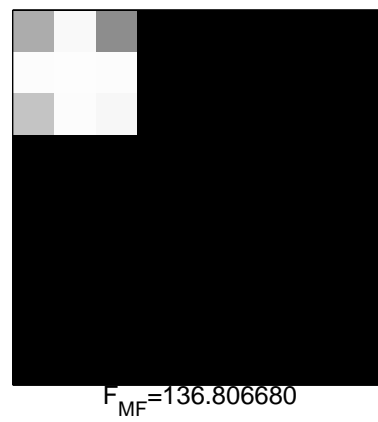
exact



loopy

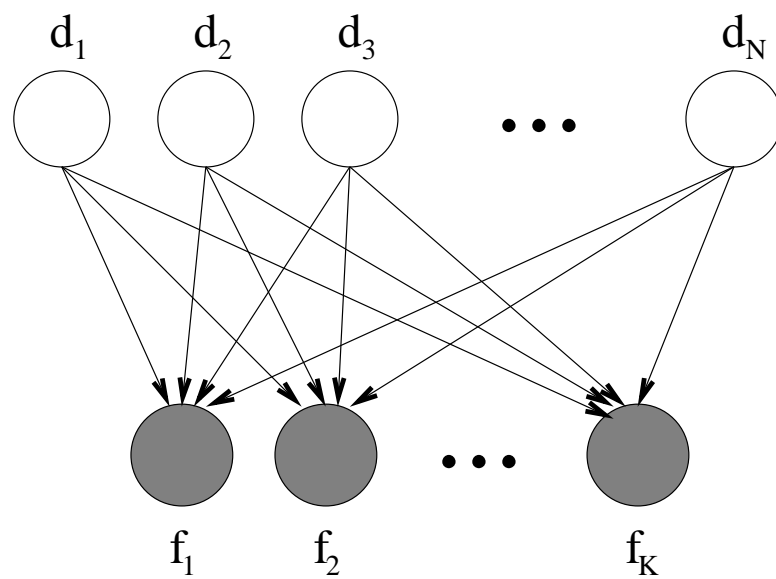


MF run 1



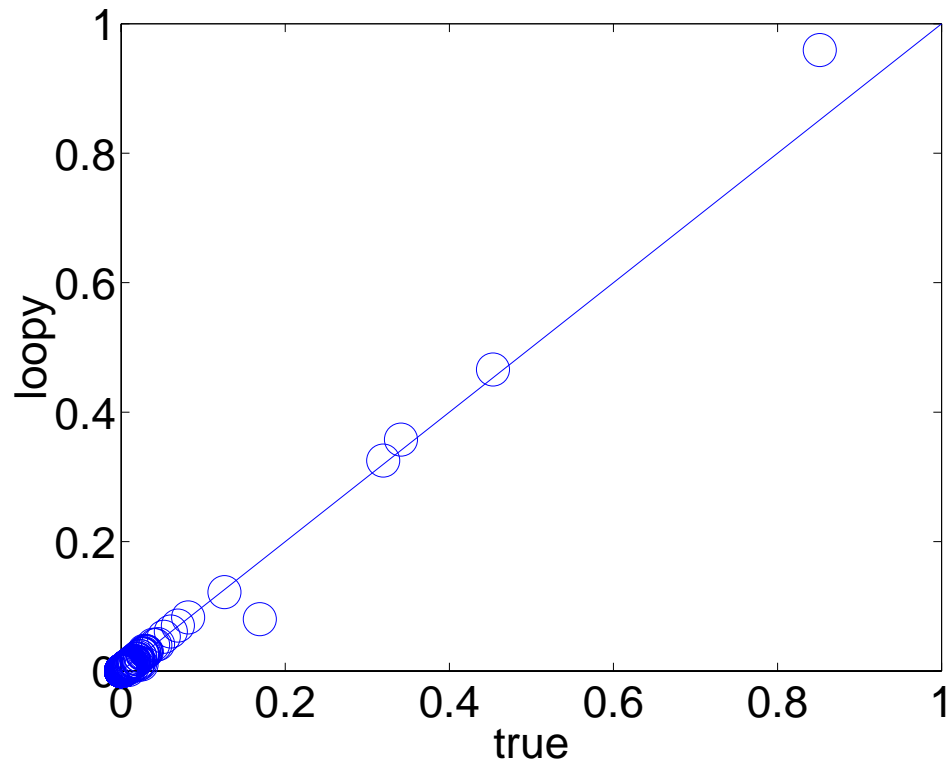
MF run 2

QMR network



approx. 600 diseases and 4000 symptoms.

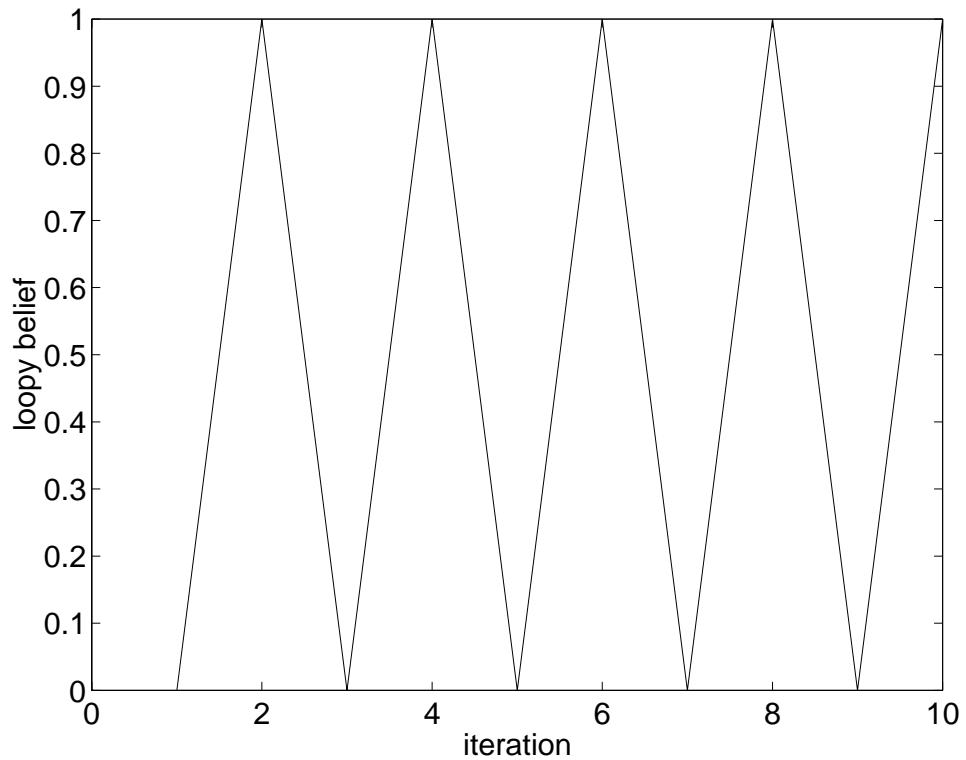
QMR results



case 32 works but case 46 does not converge.

(Murphy, Weiss and Jordan 99)

Failures of loopy BP



Two major failure modes:

- Beliefs oscillate wildly with iteration (more common).
- Beliefs converge to poor approximations.

Heuristics for avoiding oscillations

$$m_{ij}(x_j) \leftarrow \alpha \sum_{x_i} \psi_{ij}(x_i, x_j) \psi_i(x_i) \prod_{X_k \in N(X_i) \setminus X_j} m_{ki}(x_i)$$

- Asynchronous update.
- Damping.

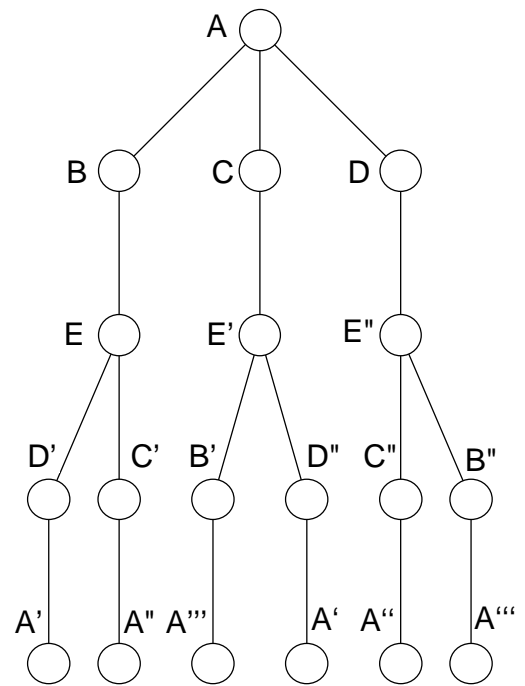
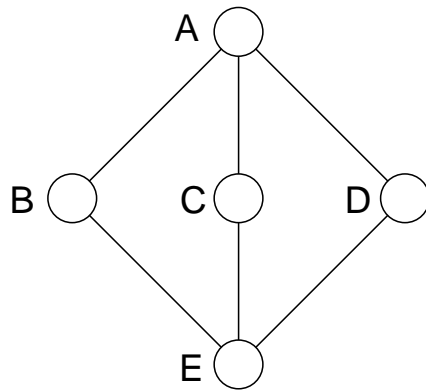
$$m_{ij} = \alpha m_{ij}^{old} + (1 - \alpha) m_{ij}^{new}$$

Analysis of loopy BP

We now understand behavior of loopy BP in:

- Single-loop graphs. Convergence proof. Correctness of decision. (Weiss 97, Aji, Horn, and McEliece 98, Forney et al 98).
- Gaussian graphs with arbitrary topology. Convergence conditions. Correctness of means. (Weiss and Freeman 99, Rusmevichientong and Van Roy 99).
- max-product in arbitrary graphs. semi-optimality of decision. (Weiss and Freeman 00)
- Average performance for LDPCs (Richardson and Urbanke 99)

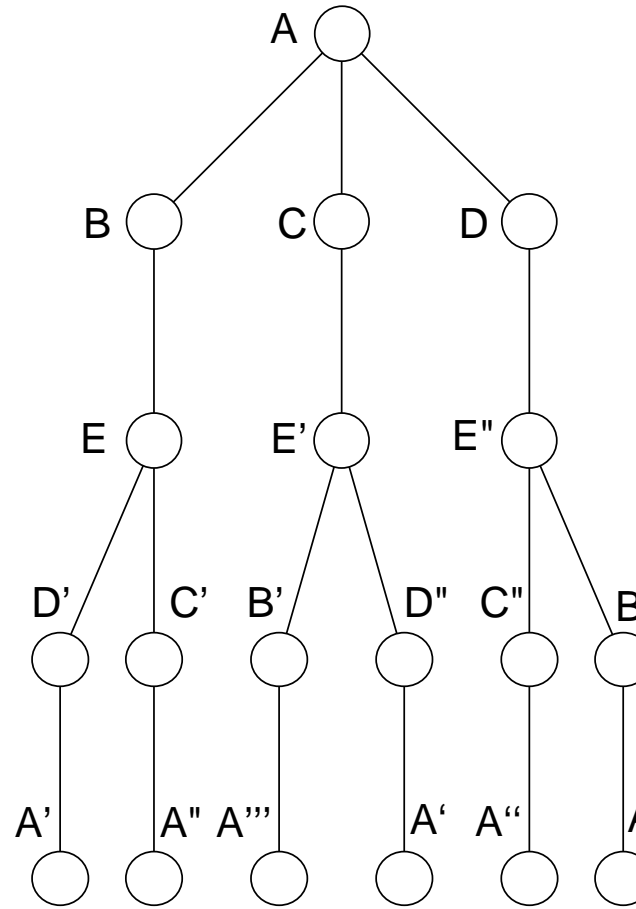
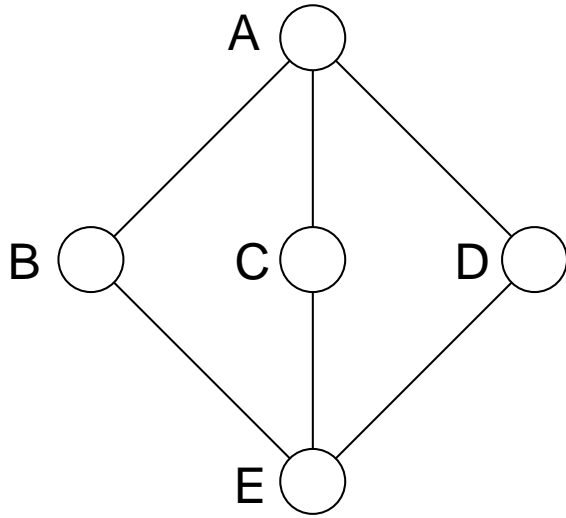
Main idea of proof:



belief propagation performs exact inference in unwrapped graph. (cf. "computation tree" Gallager 63).

unwrapped graph has same local topology as loopy graph (universal covering).

After 4 iterations



BP from a very different perspective

The BP algorithm is equivalent to an approximation introduced in statistical physics in 1935.

(Yedidia, Freeman, Weiss 01)

Variational Interpretation:

For a graphical model of *arbitrary topology and potentials* the sum-product algorithm can only converge to local stationary points of the Bethe free energy:

$$\begin{aligned} F_{\beta}(\{b_{ij}, b_i\}) = & - \sum_{ij} \sum_{x_i, x_j} b_{ij}(x_i, x_j) \ln \psi_{ij}(x_i, x_j) \\ & - \sum_i \sum_{x_i} b_i(x_i) \ln \psi_i(x_i) \\ & + \sum_{ij} \sum_{x_i, x_j} b_{ij}(x_i, x_j) \ln b_{ij}(x_i, x_j) \\ & - \sum_i (q_i - 1) \sum_{x_i} b_i(x_i) \ln b_i(x_i) \end{aligned}$$

where $b_{ij}(x_i, x_j)$ is joint belief of x_i, x_j , q_i is degree of node i .

Yedidia J.S., Freeman W.T. and Weiss Y.

“Generalized Belief Propagation” NIPS 2000.

What does Bethe teach us about BP?

- Existence of fixed-point for arbitrary graph.
- Conditions for uniqueness of fixed-points. Phase transitions.
- goodness of approximation ???

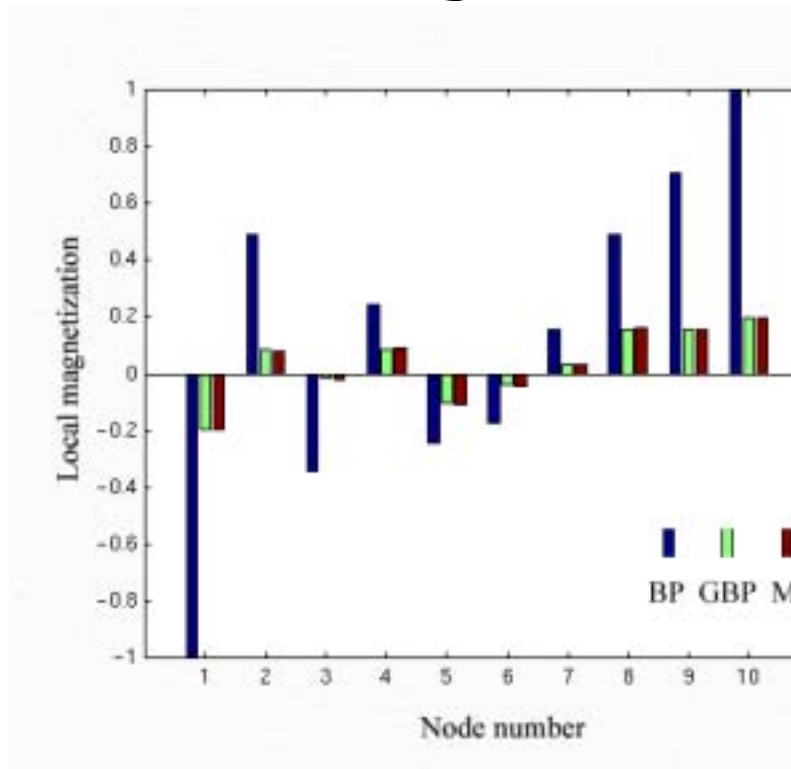
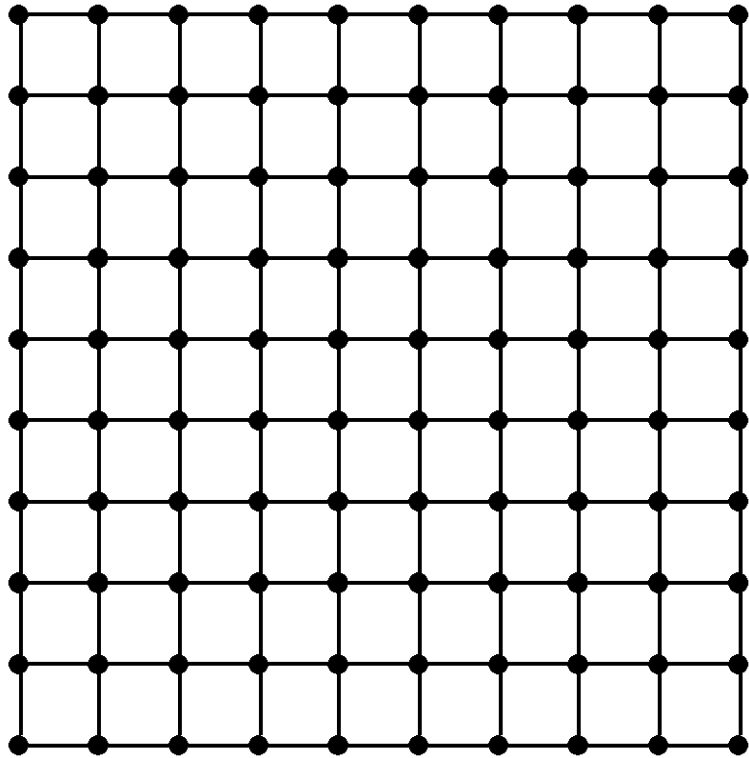
Beyond BP

- Direct minimization of Bethe free energy. Guaranteed to converge. (Yuille 01, Teh and Welling 01, Kappen et al. 02).
- “Convexified” free energy. (Wainwright et al. 02) Minor modification on max-product BP. Guaranteed to find **global** optimum if it converges.
- Generalized belief propagation.

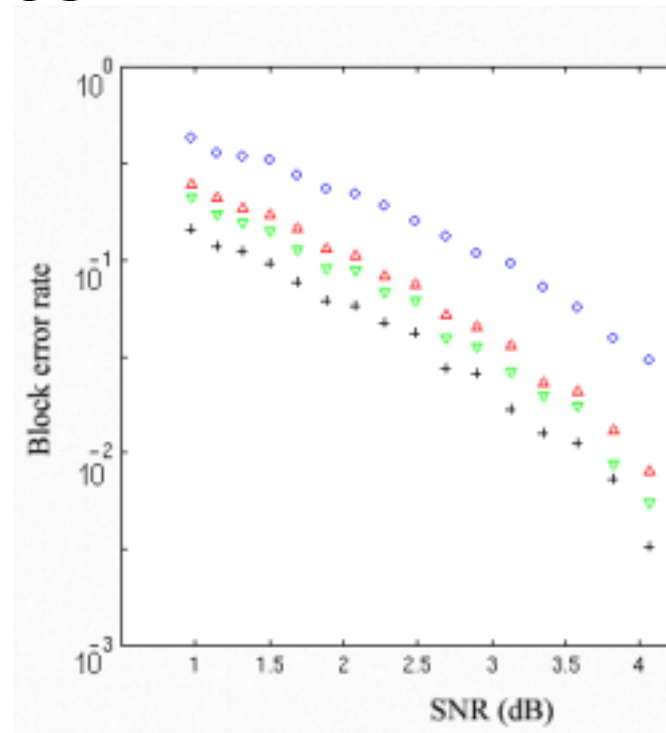
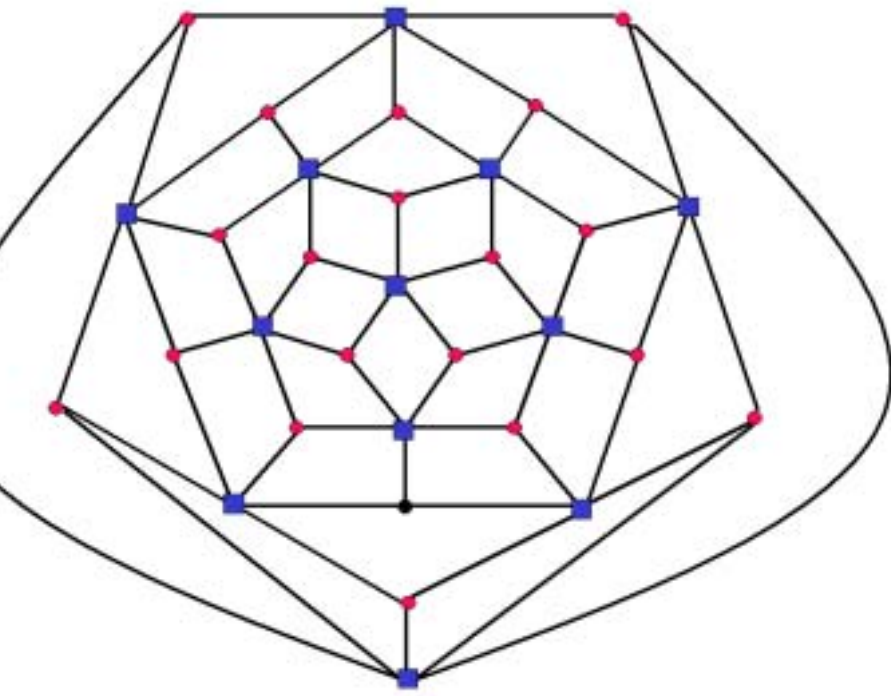
Generalized Belief Propagation

- A message passing scheme for minimizing Kikuchi approximations.
- messages are functions of clusters of variables, not single variables.
- message update includes sums, products and **ratios**
- several algorithms depending on which Lagrange multipliers are defined.

Comparison on spin glass



Comparison on error-correcting codes



Conclusions

- The need for approximate inference.
- loopy belief propagation. Was not supposed to work but dramatic empirical results — turbo codes.
- Theory of loopy BP: unwrapped tree and Bethe free energy.
- Kikuchi free energies extend Bethe. Can be minimized by generalized BP algorithms.
- Simple, universally applicable approximate inference algorithms.