

Machine Learning for Data Science (CS4786)

Lecture 20

Finish HMM, Inference in Graphical Models

Course Webpage :

<http://www.cs.cornell.edu/Courses/cs4786/2016fa/>

ANNOUNCEMENT

- No lecture Tuesday, Nov 8th!
- Next Thursday Nov 10th, guest lecture by Prof. Kilian Weinberger on TSNE

BAYESIAN NETWORKS

- Directed acyclic graph (DAG): $G = (V, E)$
- Joint distribution P_θ over X_1, \dots, X_n that factorizes over G :

$$P_\theta(X_1, \dots, X_n) = \prod_{i=1}^n P_\theta(X_i | \text{Parent}(X_i))$$

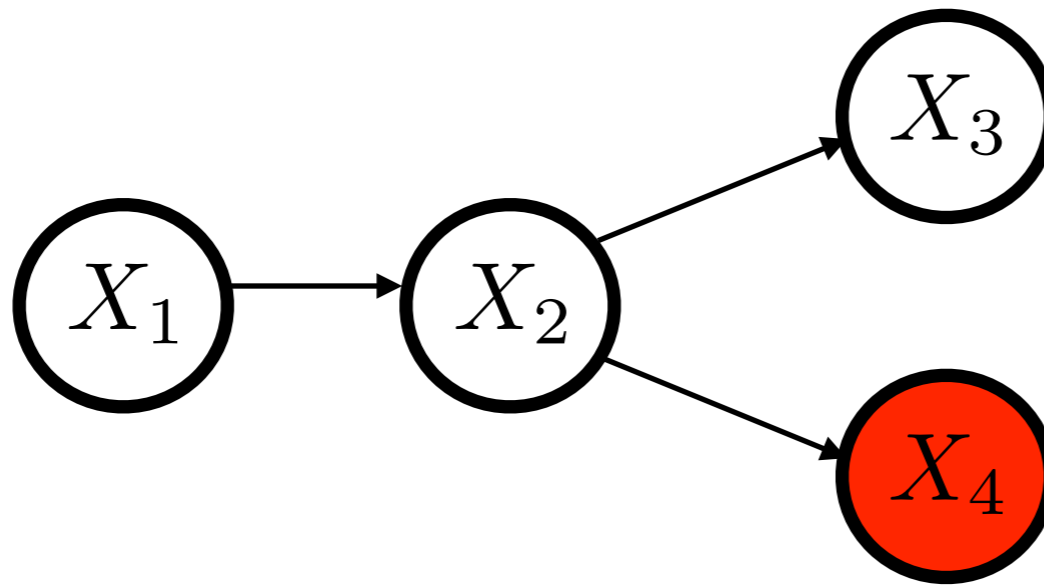
- Hence Bayesian Networks are specified by G along with CPD's over the variables (given their parents)

VARIABLE ELIMINATION: EXAMPLES

- Marginals are enough:

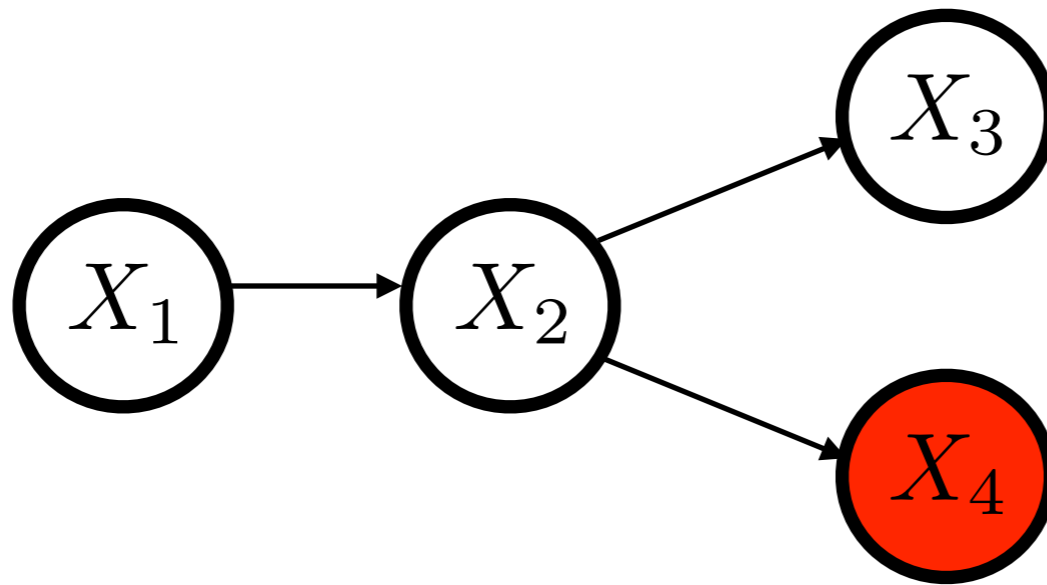
$$P(X_j = x_j, X_k = x_k | X_i = x_i, X_h = x_h) = \frac{P(X_j = x_j, X_k = x_k, X_i = x_i, X_h = x_h)}{P(X_i = x_i, X_h = x_h)}$$

VARIABLE ELIMINATION: EXAMPLES



$P(\text{Given variables}) = \text{Sum over all other variables } (P(\text{All variables}))$
 $= \text{Sum over all other variables } (\text{Product } P(X_i | \text{Parents}(X_i)))$

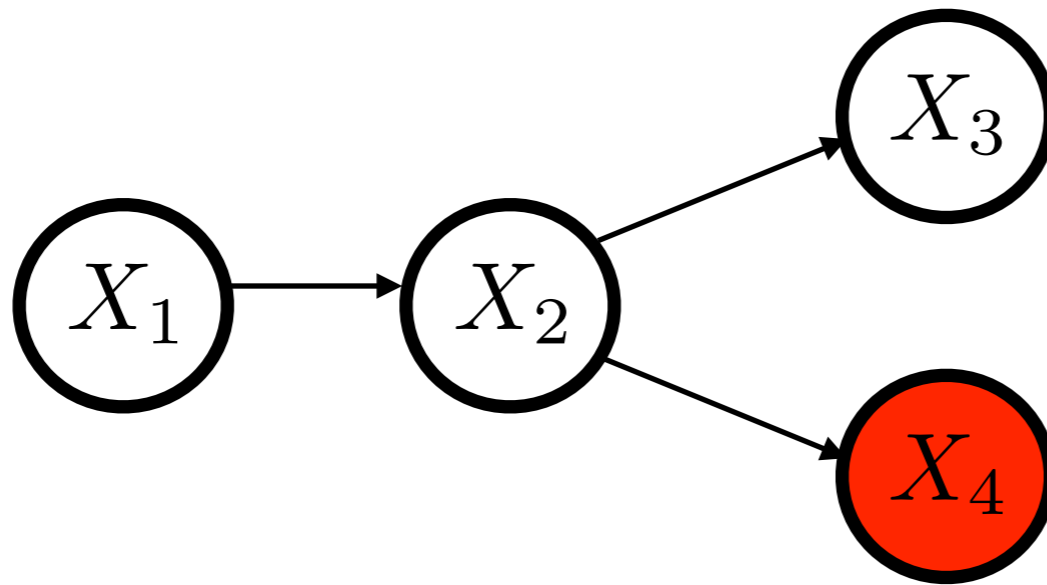
VARIABLE ELIMINATION: EXAMPLES



$P(\text{Given variables}) = \text{Sum over all other variables } (P(\text{All variables}))$
 $= \text{Sum over all other variables } (\text{Product } P(X_i | \text{Parents}(X_i)))$

$$P(X_4) = \sum_{x_1} \sum_{x_2} \sum_{x_3} P(X_1 = x_1, X_2 = x_2, X_3 = x_3, X_4)$$

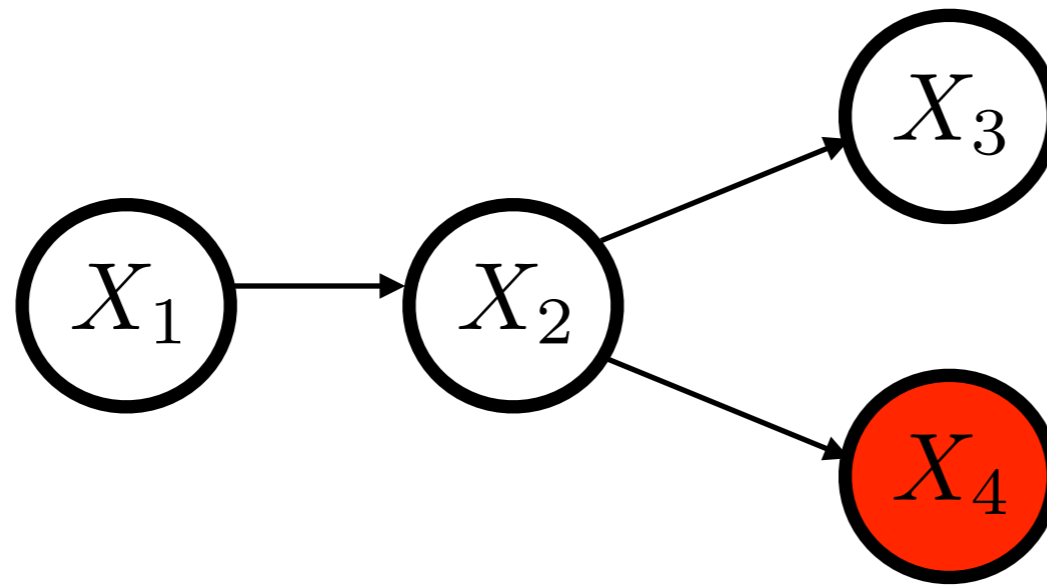
VARIABLE ELIMINATION: EXAMPLES



$P(\text{Given variables}) = \text{Sum over all other variables } (P(\text{All variables}))$
 $= \text{Sum over all other variables } (\text{Product } P(X_i | \text{Parents}(X_i)))$

$$\begin{aligned} P(X_4) &= \sum_{x_1} \sum_{x_2} \sum_{x_3} P(X_1 = x_1, X_2 = x_2, X_3 = x_3, X_4) \\ &= \sum_{x_1} \sum_{x_2} \sum_{x_3} (P(X_1 = x_1) \cdot P(X_2 = x_2 | X_1 = x_1) \cdot P(X_3 = x_3 | X_2 = x_2) \cdot P(X_4 | X_2 = x_2)) \end{aligned}$$

VARIABLE ELIMINATION: EXAMPLES



$P(\text{Given variables}) = \text{Sum over all other variables } (P(\text{All variables}))$
 $= \text{Sum over all other variables } (\text{Product } P(X_i | \text{Parents}(X_i)))$

$$\begin{aligned} P(X_4) &= \sum_{x_1} \sum_{x_2} \sum_{x_3} P(X_1 = x_1, X_2 = x_2, X_3 = x_3, X_4) \\ &= \sum_{x_1} \sum_{x_2} \sum_{x_3} (P(X_1 = x_1) \cdot P(X_2 = x_2 | X_1 = x_1) \cdot P(X_3 = x_3 | X_2 = x_2) \cdot P(X_4 | X_2 = x_2)) \\ &= \sum_{x_1} \left(P(X_1 = x_1) \sum_{x_2} \left(P(X_2 = x_2 | X_1 = x_1) P(X_4 | X_2 = x_2) \left(\sum_{x_3} P(X_3 = x_3 | X_2 = x_2) \right) \right) \right) \end{aligned}$$

VARIABLE ELIMINATION: BAYESIAN NETWORK

Initialize **List** with conditional probability distributions

Pick an order of elimination I for remaining variables

For each $X_i \in I$

 Find distributions in **List** containing variable X_i and remove them

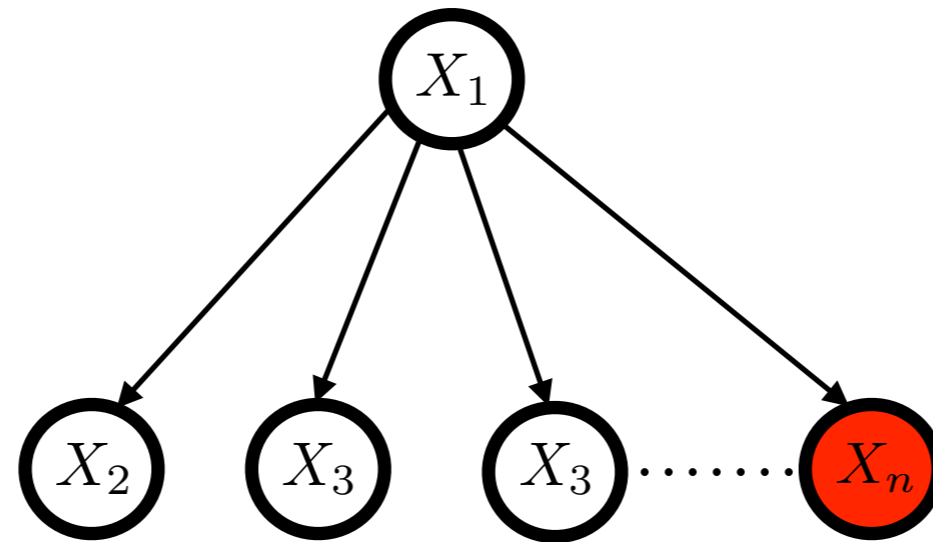
 Define new distribution as the sum (over values of X_i) of the product of these distributions

 Place the new distribution on **List**

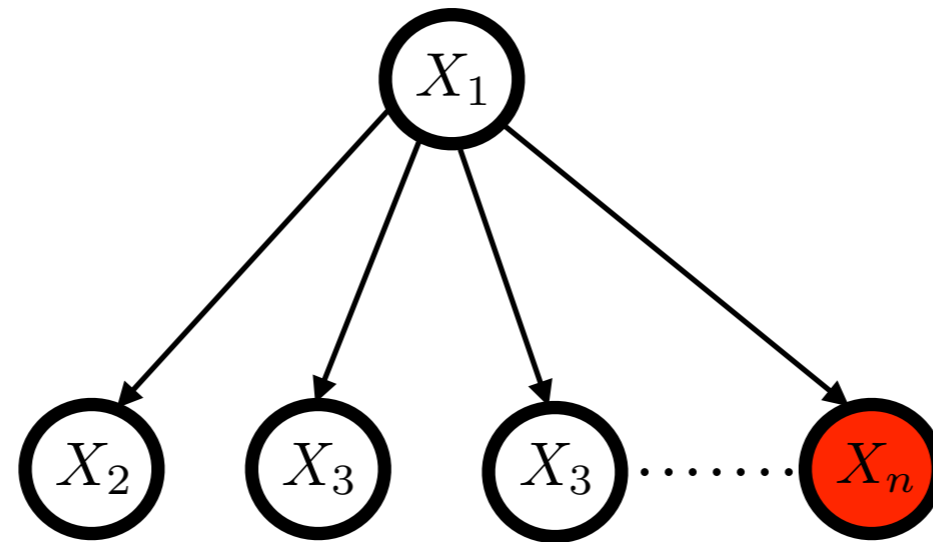
End

Return **List**

VARIABLE ELIMINATION: ORDER MATTERS

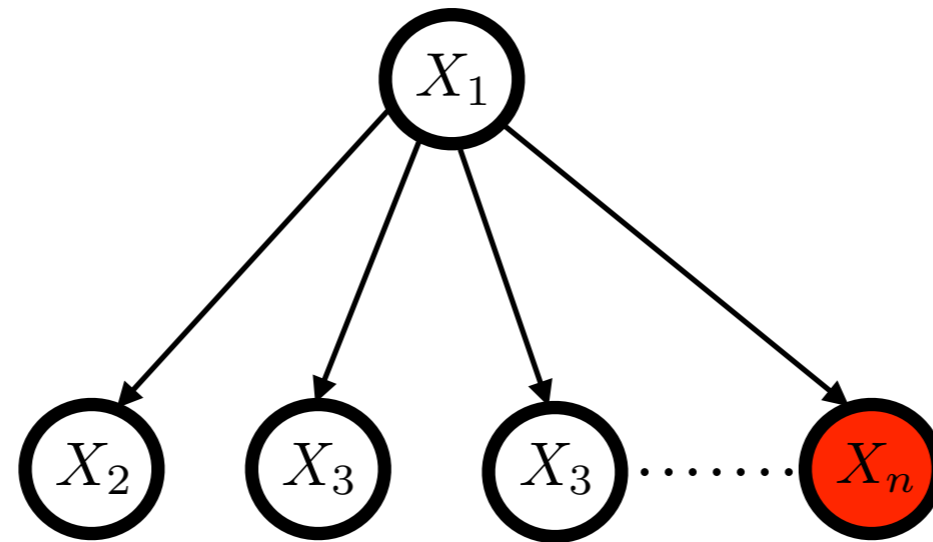


VARIABLE ELIMINATION: ORDER MATTERS



List initialized to : $\{P(X_1), P(X_2|X_1), P(X_3|X_1), \dots, P(X_n|X_1)\}$

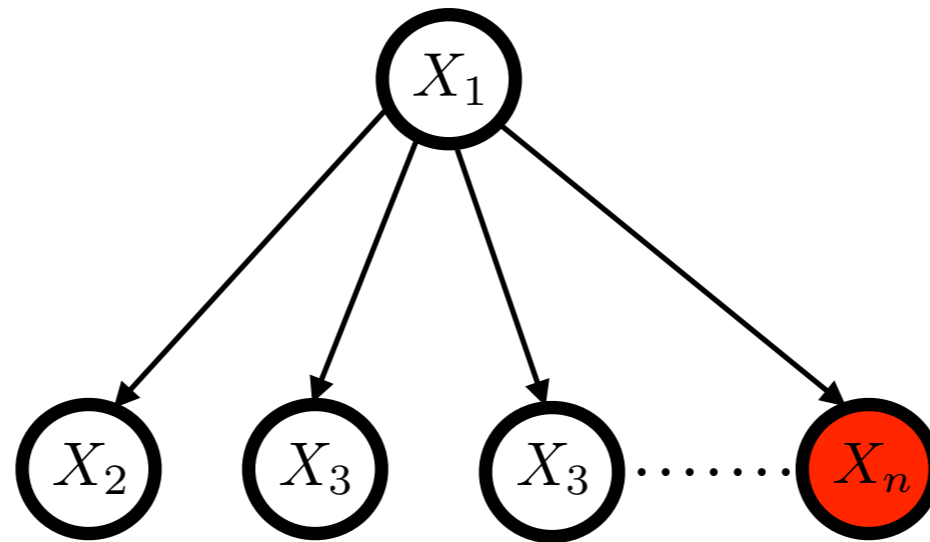
VARIABLE ELIMINATION: ORDER MATTERS



List initialized to : $\{P(X_1), P(X_2|X_1), P(X_3|X_1), \dots, P(X_n|X_1)\}$

Say $I = (1, 2, 3, \dots, n-1)$

VARIABLE ELIMINATION: ORDER MATTERS

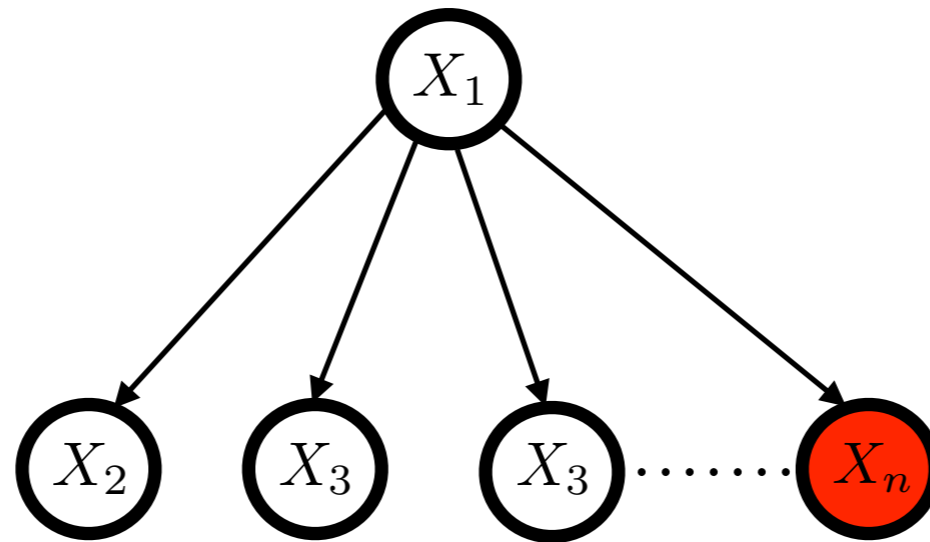


List initialized to : $\{P(X_1), P(X_2|X_1), P(X_3|X_1), \dots, P(X_n|X_1)\}$

Say $I = (1, 2, 3, \dots, n-1)$

Iteration 1: Eliminate X_1

VARIABLE ELIMINATION: ORDER MATTERS



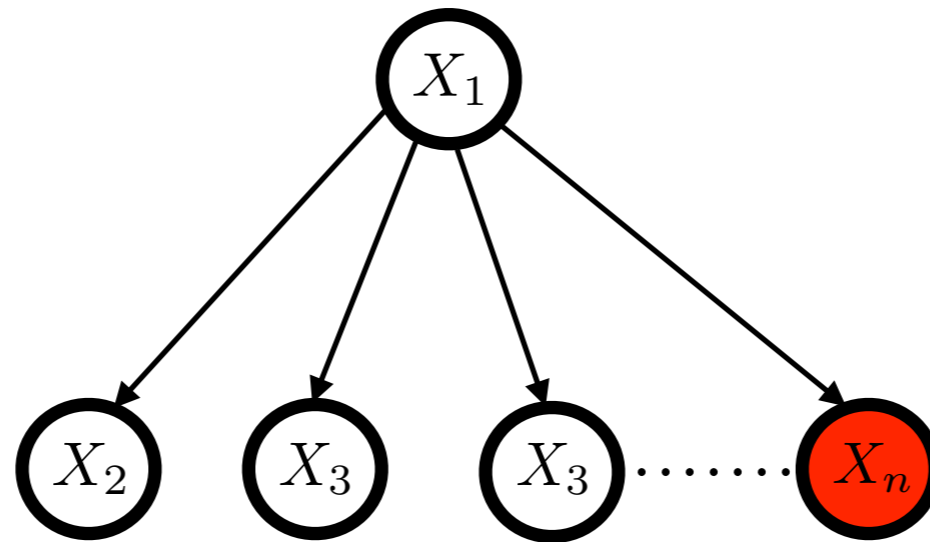
List initialized to : $\{P(X_1), P(X_2|X_1), P(X_3|X_1), \dots, P(X_n|X_1)\}$

Say $I = (1, 2, 3, \dots, n-1)$

Iteration 1: Eliminate X_1

All terms in list involve X_1 so remove all of them

VARIABLE ELIMINATION: ORDER MATTERS



List initialized to : $\{P(X_1), P(X_2|X_1), P(X_3|X_1), \dots, P(X_n|X_1)\}$

Say $I = (1, 2, 3, \dots, n-1)$

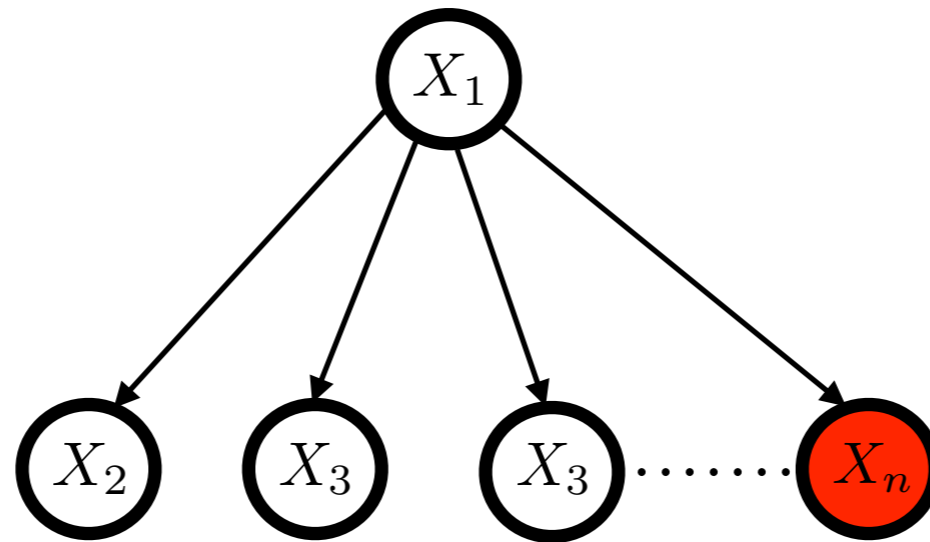
Iteration 1: Eliminate X_1

All terms in list involve X_1 so remove all of them

Replace them by table:

$$L_1(x_1, \dots, x_n) = \sum_{x_1} \left(P(X_1 = x_1) \prod_{t=2}^n P(X_t = x_t | X_1 = x_1) \right)$$

VARIABLE ELIMINATION: ORDER MATTERS



List initialized to $\{P(Y_1), P(Y_2|Y_1), P(Y_3|Y_1), \dots, P(Y_n|X_1)\}$

Say

What is the size of table L1?

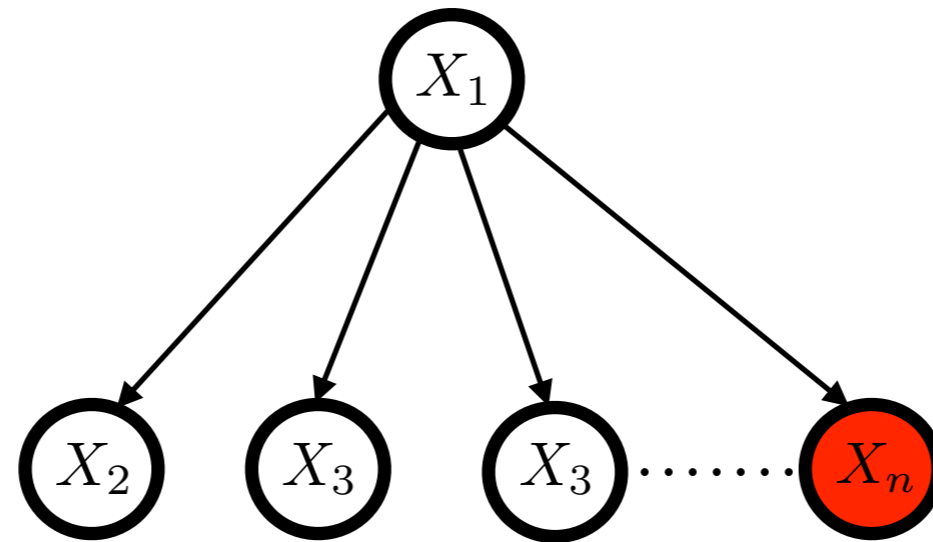
Iteration 1: Eliminate X_1

All terms in list involve X_1 so remove all of them

Replace them by table:

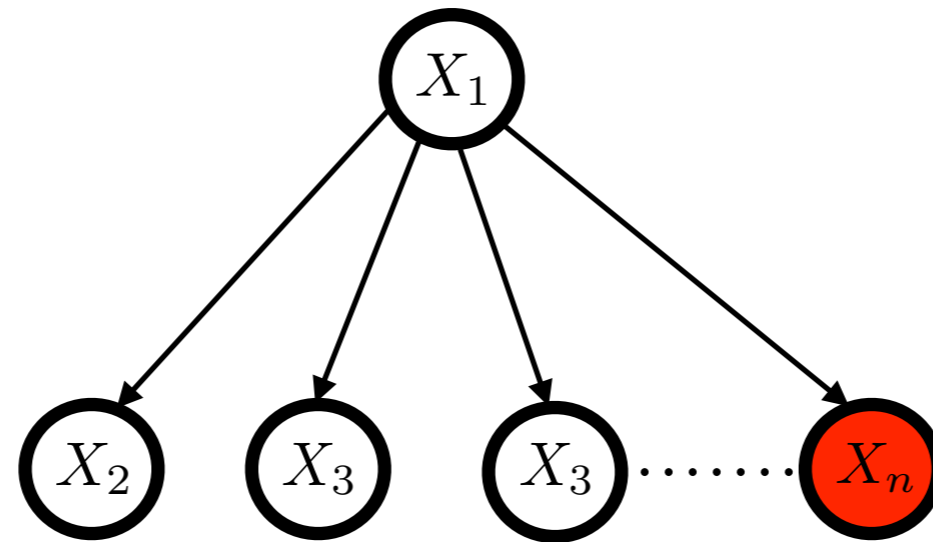
$$L_1(x_1, \dots, x_n) = \sum_{x_1} \left(P(X_1 = x_1) \prod_{t=2}^n P(X_t = x_t | X_1 = x_1) \right)$$

VARIABLE ELIMINATION: ORDER MATTERS



List initialized to : $\{P(X_1), P(X_2|X_1), P(X_3|X_1), \dots, P(X_n|X_1)\}$

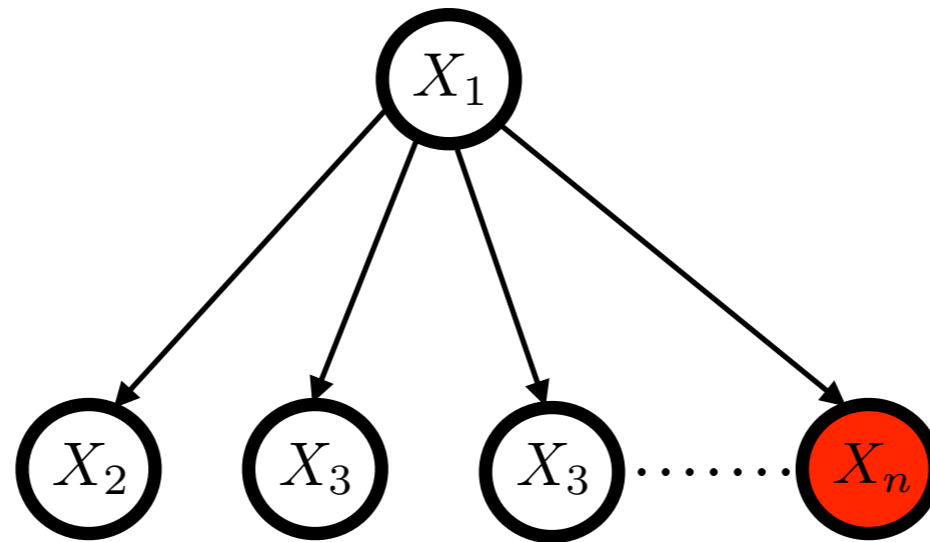
VARIABLE ELIMINATION: ORDER MATTERS



List initialized to : $\{P(X_1), P(X_2|X_1), P(X_3|X_1), \dots, P(X_n|X_1)\}$

Say $I = (n-1, n-2, \dots, 1)$

VARIABLE ELIMINATION: ORDER MATTERS

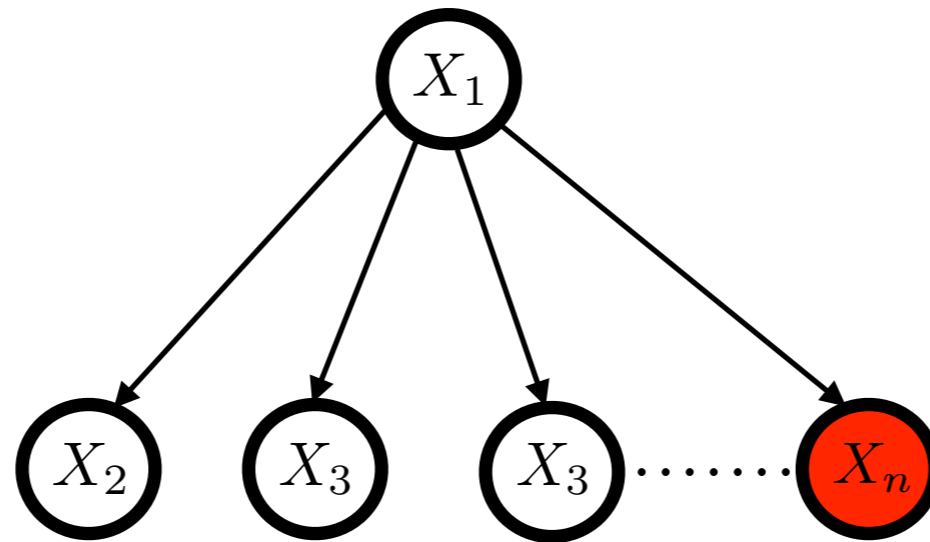


List initialized to : $\{P(X_1), P(X_2|X_1), P(X_3|X_1), \dots, P(X_n|X_1)\}$

Say $I = (n-1, n-2, \dots, 1)$

Iteration 1: Eliminate X_{n-1}

VARIABLE ELIMINATION: ORDER MATTERS



List initialized to : $\{P(X_1), P(X_2|X_1), P(X_3|X_1), \dots, P(X_n|X_1)\}$

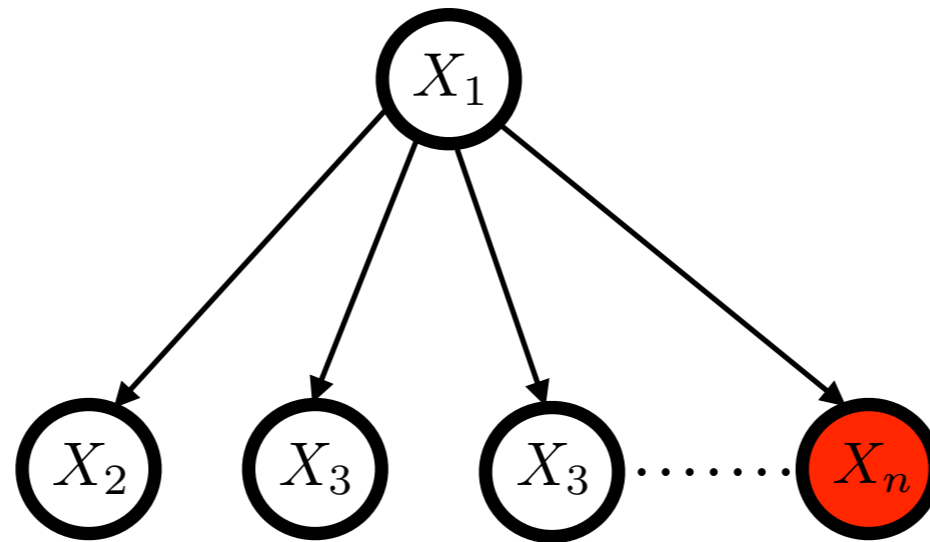
Say $I = (n-1, n-2, \dots, 1)$

Iteration 1: Eliminate X_{n-1}

Remove $P(X_n|X_1)$ from List and replace by

$$L_{n-1}(x_1) = \sum_{x_{n-1}} P(X_{n-1}|X_1 = x_1) = 1$$

VARIABLE ELIMINATION: ORDER MATTERS



List initialized to : $\{P(X_1), P(X_2|X_1), P(X_3|X_1), \dots, P(X_n|X_1)\}$

Say $I = (n-1, n-2, \dots, 1)$

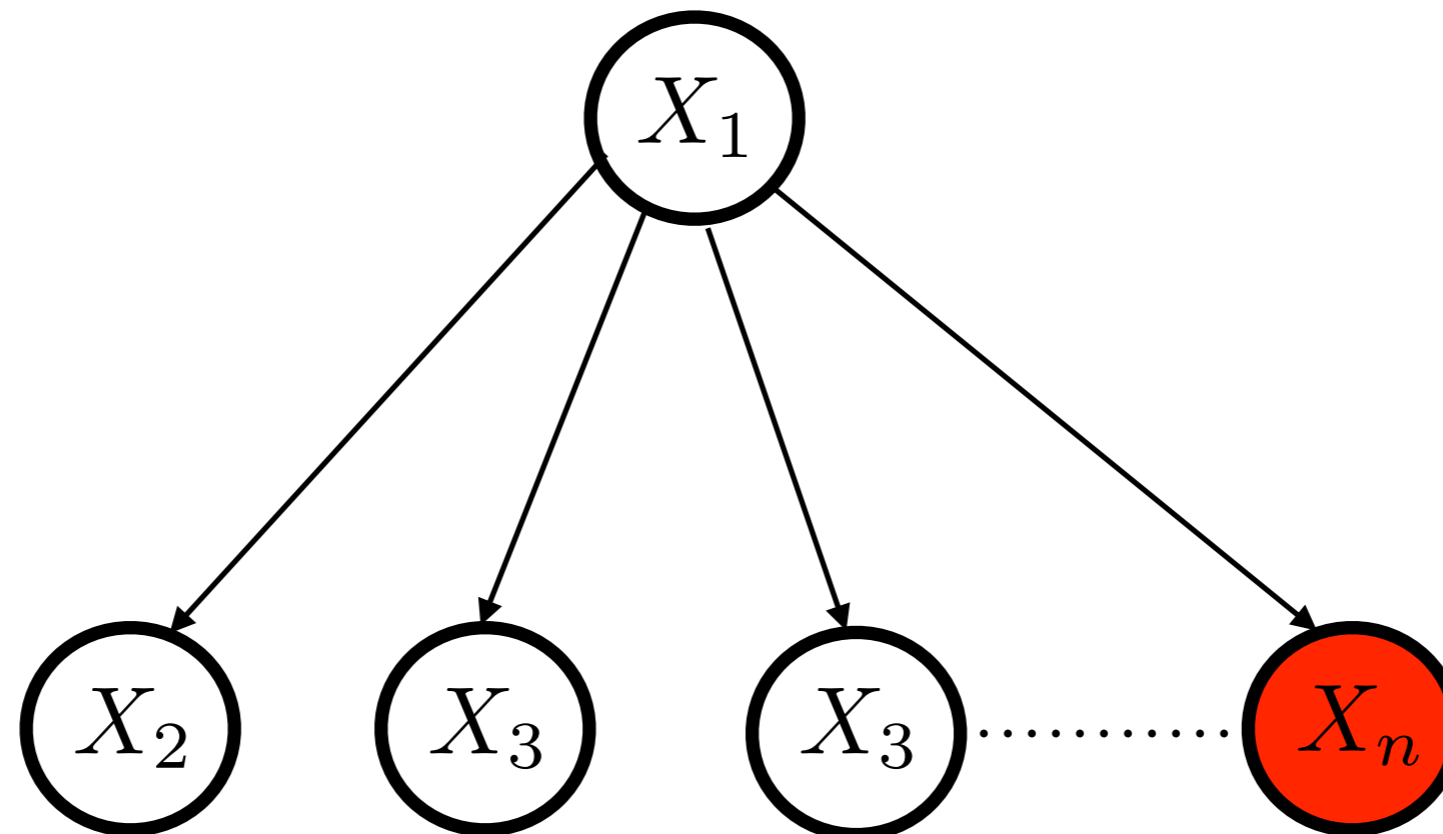
Iteration 1: Eliminate X_{n-1}

Remove $P(X_n|X_1)$ from List and replace by

$$L_{n-1}(x_1) = \sum_{x_{n-1}} P(X_{n-1}|X_1 = x_1) = 1$$

All the way up to X_2 we replace by all ones message
In then end we only have $P(X_1), P(X_n|X_1)$

VARIABLE ELIMINATION: ORDER MATTERS



Right order: $O(n)$

Wrong order: $O(2^n)$

MESSAGE PASSING

- Often we need more than one marginal computation
- Can we exploit structure and compute these intermediate terms that can be reused?

Eg. forward backward algorithm for HMM

BELIEF PROPAGATION

- Think of vertices in graphical model as nodes in a network
- On every round, each node sends and receives messages from all its neighbors

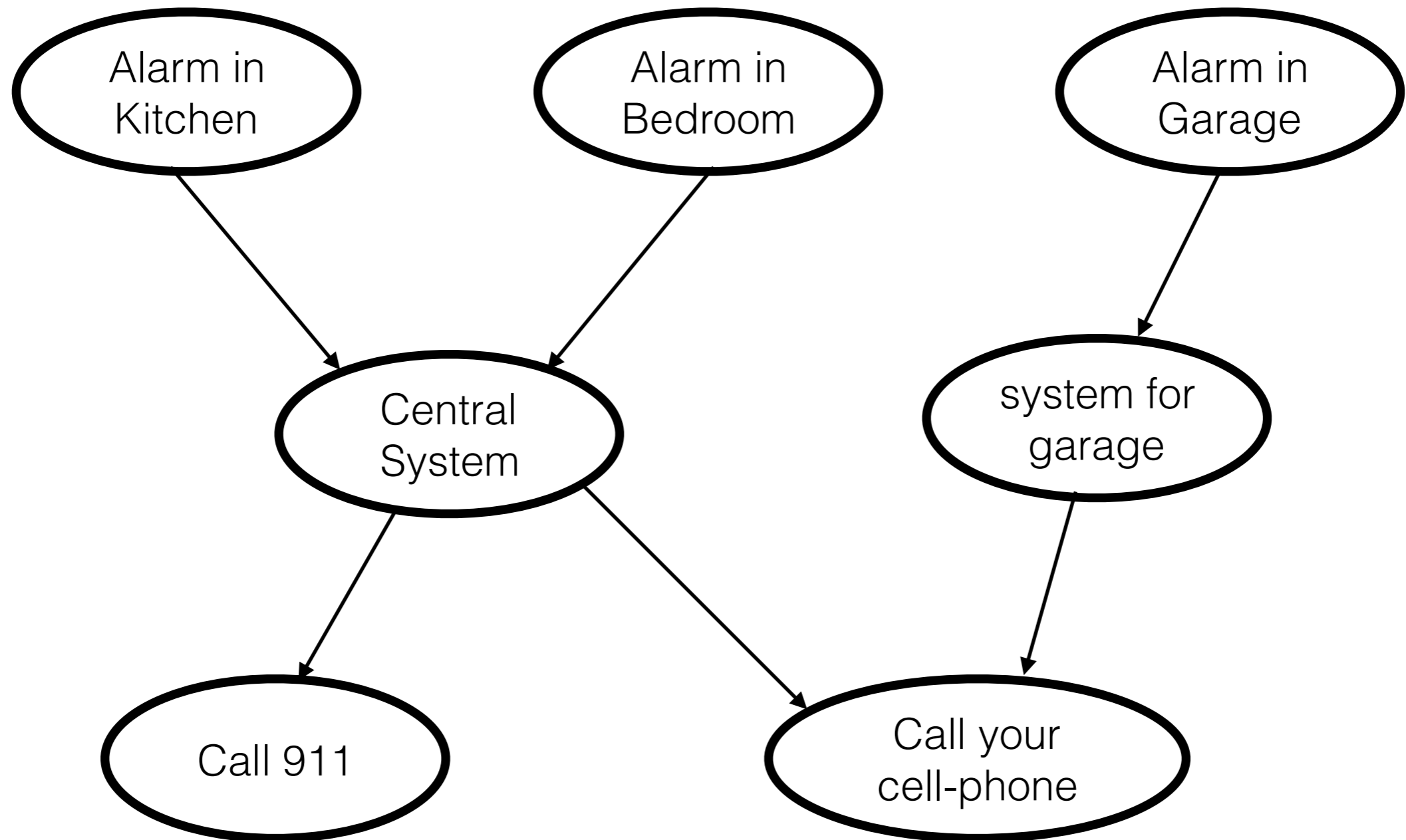
BELIEF PROPAGATION

- This is specific to Bayesian networks
 - Messages to parents: belief about parent's value
 - Messages to Children: belief about self

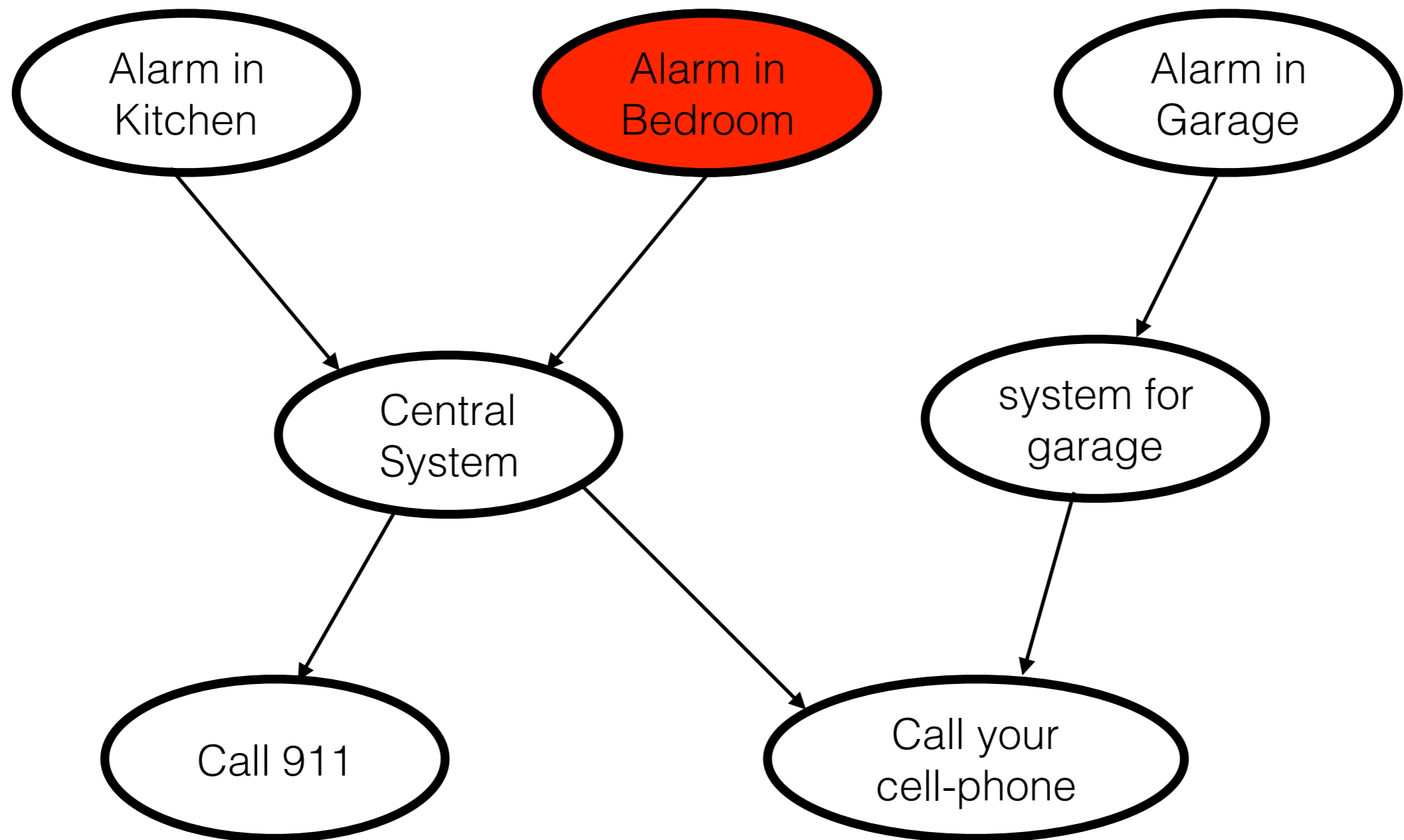
BELIEF PROPAGATION

- Evidence E_i for node X_i can be seen as a priori belief of each node about itself.
 - if unobserved, set it to uniform distribution
 - if $X_i = x_i$ is observed,
set $E_i(x_i) = 1$ and $E_i(x) = 0$ for any other x

BELIEF PROPAGATION

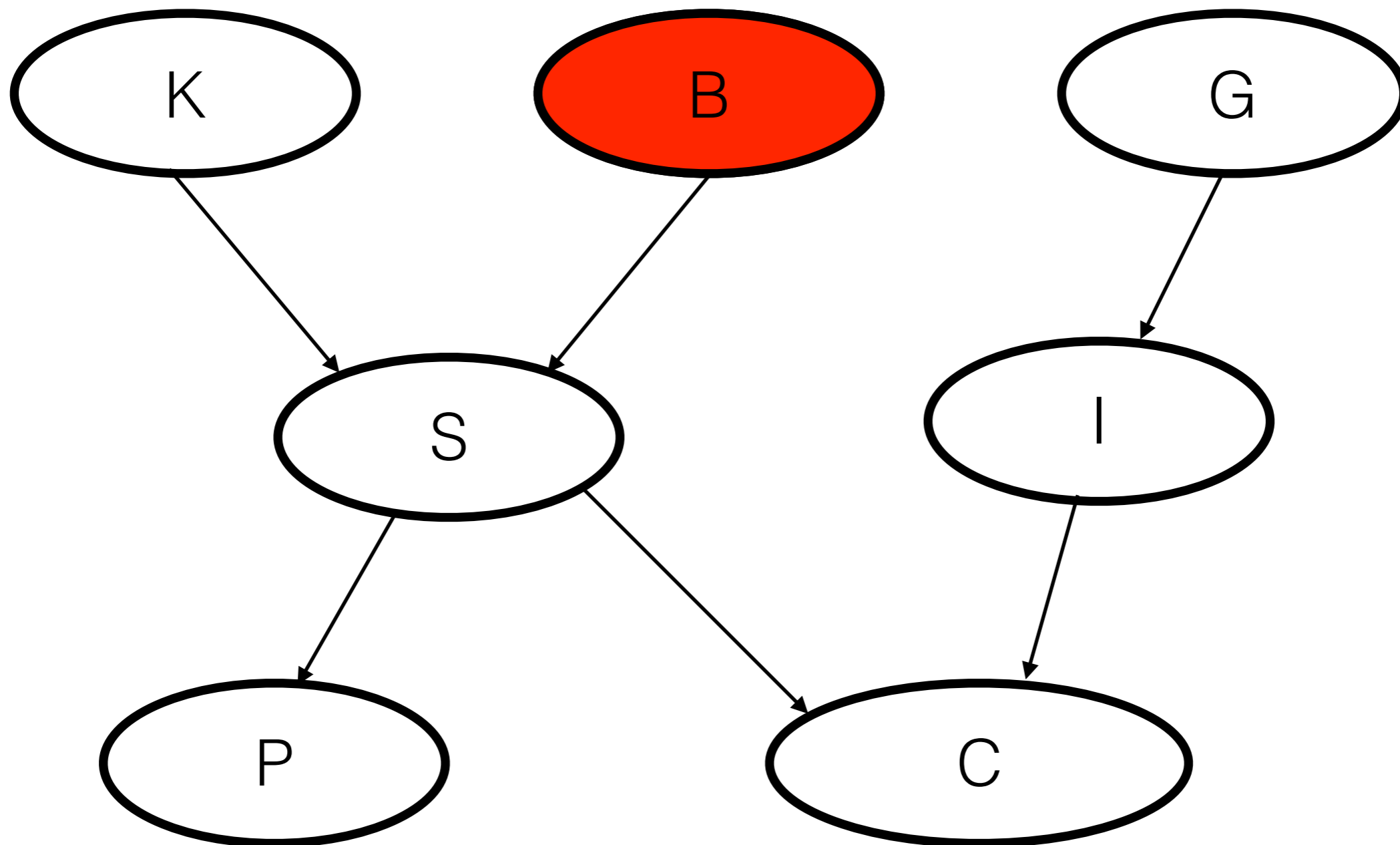


BELIEF PROPAGATION

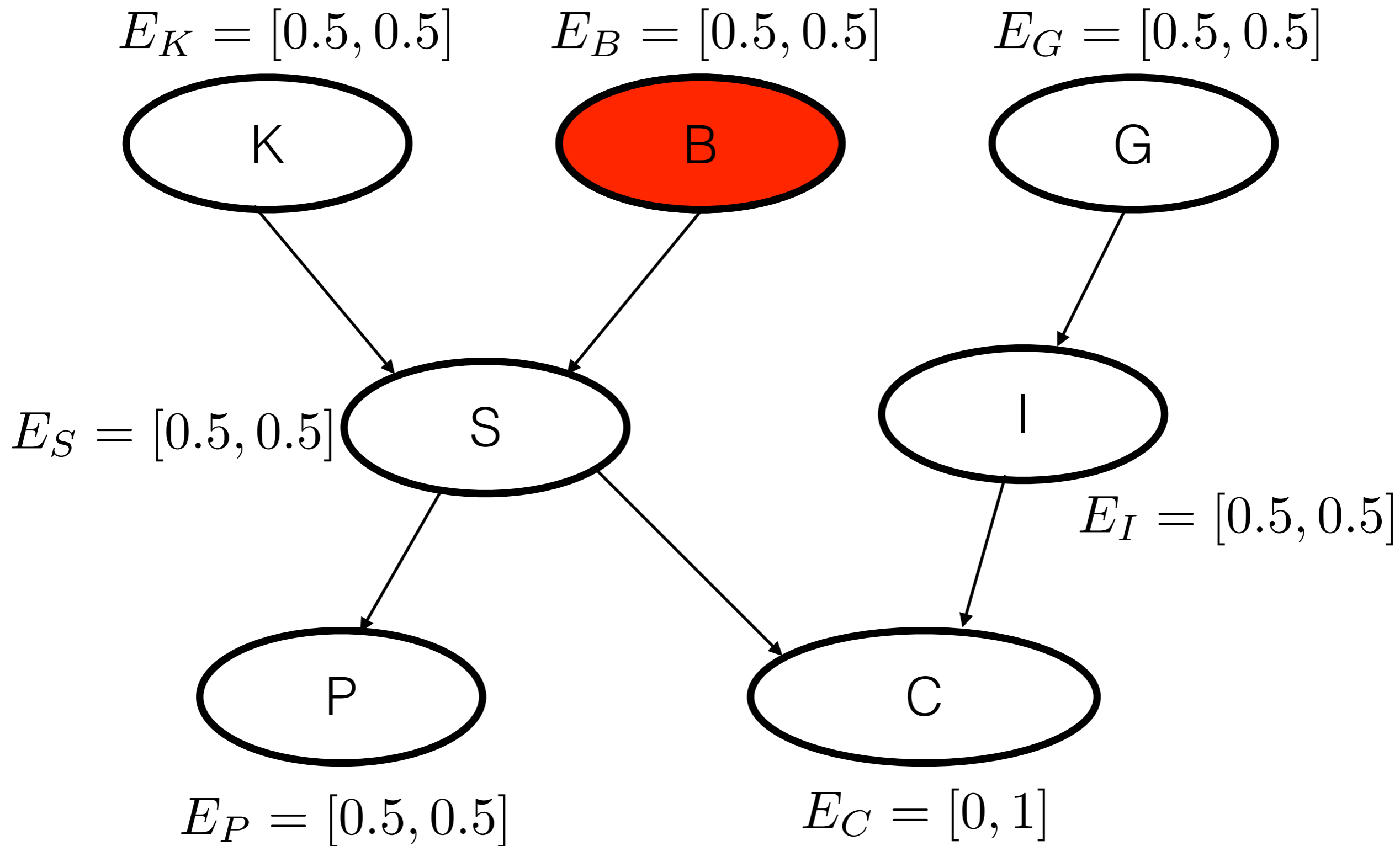


You receive phone call

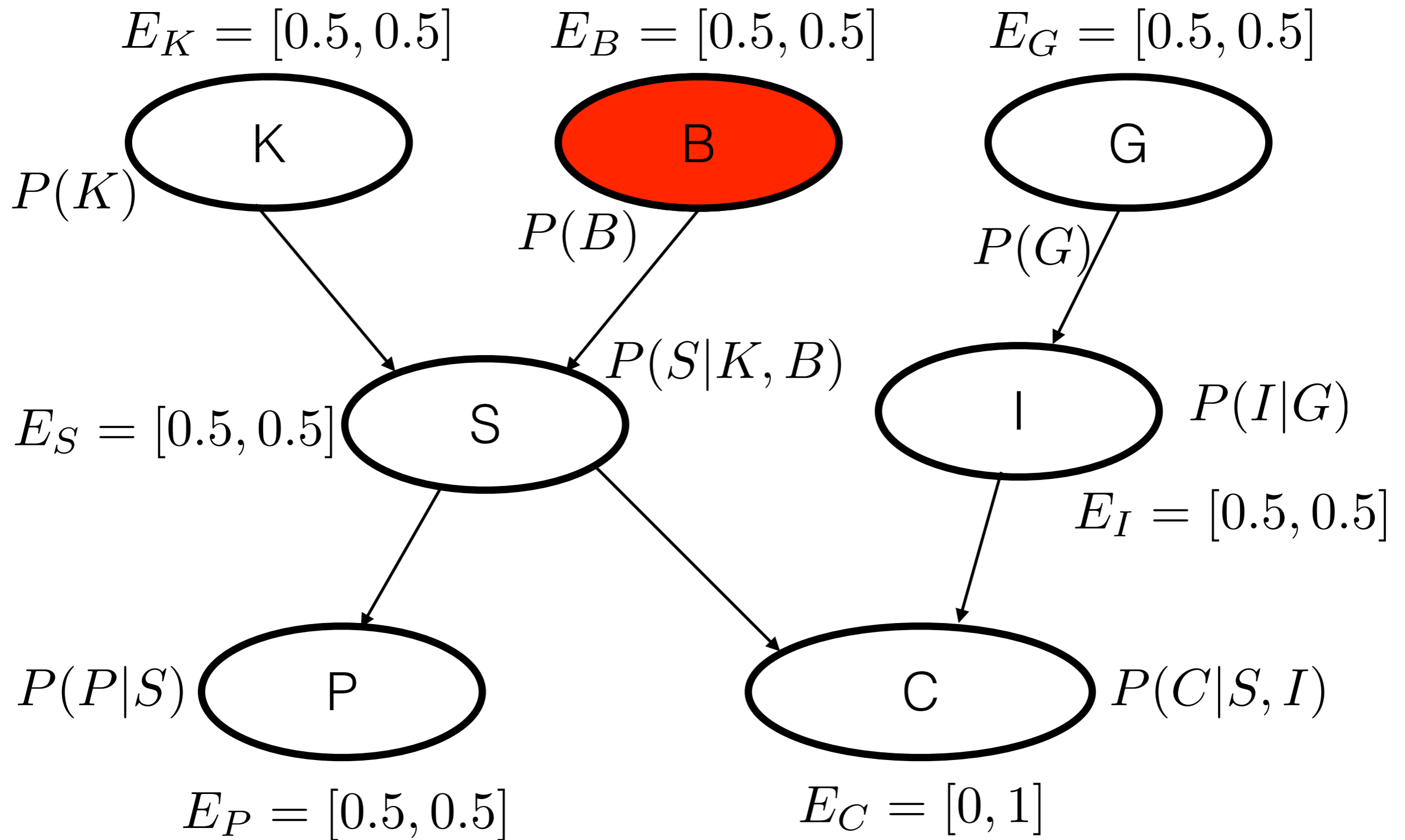
BELIEF PROPAGATION



BELIEF PROPAGATION



BELIEF PROPAGATION



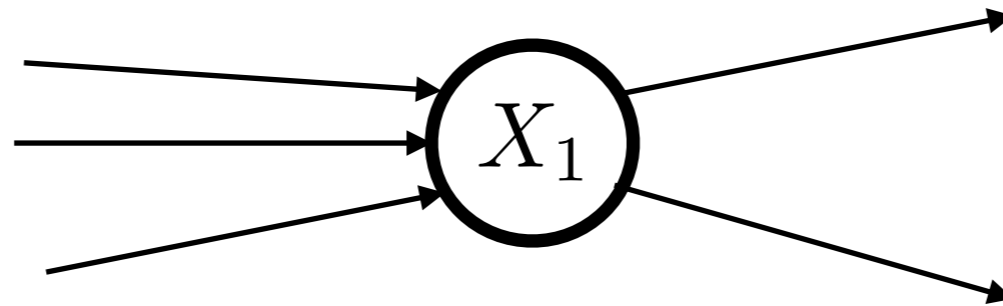
BELIEF PROPAGATION

- On each round:
Use evidence + messages from previous round to compute new message

BELIEF PROPAGATION

- On each round:

Use evidence + messages from previous round to compute new message

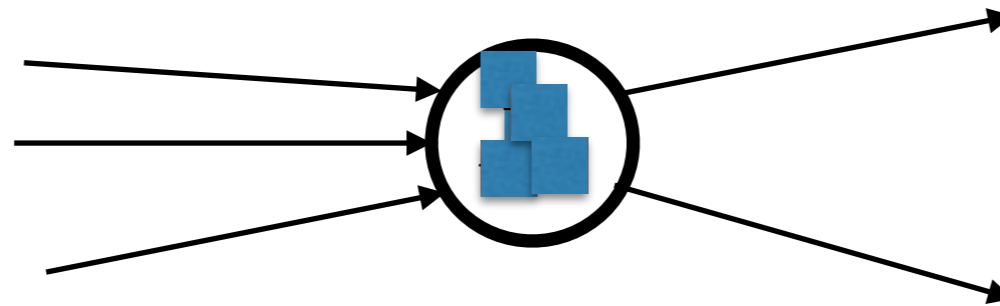


BELIEF PROPAGATION

- On each round:

Use evidence + messages from previous round to compute new message

Round $i-1$

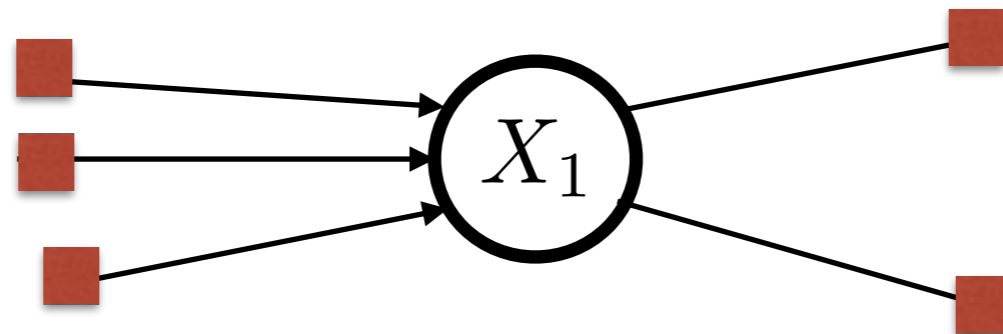


BELIEF PROPAGATION

- On each round:

Use evidence + messages from previous round to compute new message

Round i



BELIEF PROPAGATION

BELIEF PROPAGATION

Message from node X_i to **Parent** X_j on round t

$$M_{i \rightarrow j}^t(x_j) = \sum_{x_i, \text{Parents}(X_i) \setminus X_j} E_{X_i}(x_i) P(X_i = x_i | \text{Parents}(X_i))$$

(product of all messages but one from X_j)
from previous round (t-1)

BELIEF PROPAGATION

Message from node X_i to **Parent** X_j on round t

$$M_{i \rightarrow j}^t(x_j) = \sum_{x_i, \text{Parents}(X_i) \setminus X_j} E_{X_i}(x_i) P(X_i = x_i | \text{Parents}(X_i)) \text{ (product of all messages but one from } X_j \text{)} \\ \text{from previous round (t-1)}$$

Message from node X_i to **Child** X_k on round t

$$M_{i \rightarrow k}^t(x_i) = \sum_{\text{Parents}(X_i)} E_{X_i}(x_i) P(X_i = x_i | \text{Parents}(X_i)) \text{ (product of all messages but one from } X_j \text{)} \\ \text{from previous round (t-1)}$$

BELIEF PROPAGATION

Message from node X_i to **Parent** X_j on round t

$$M_{i \rightarrow j}^t(x_j) = \sum_{x_i, \text{Parents}(X_i) \setminus X_j} E_{X_i}(x_i) P(X_i = x_i | \text{Parents}(X_i)) \text{ (product of all messages but one from } X_j \text{)} \\ \text{from previous round (t-1)}$$

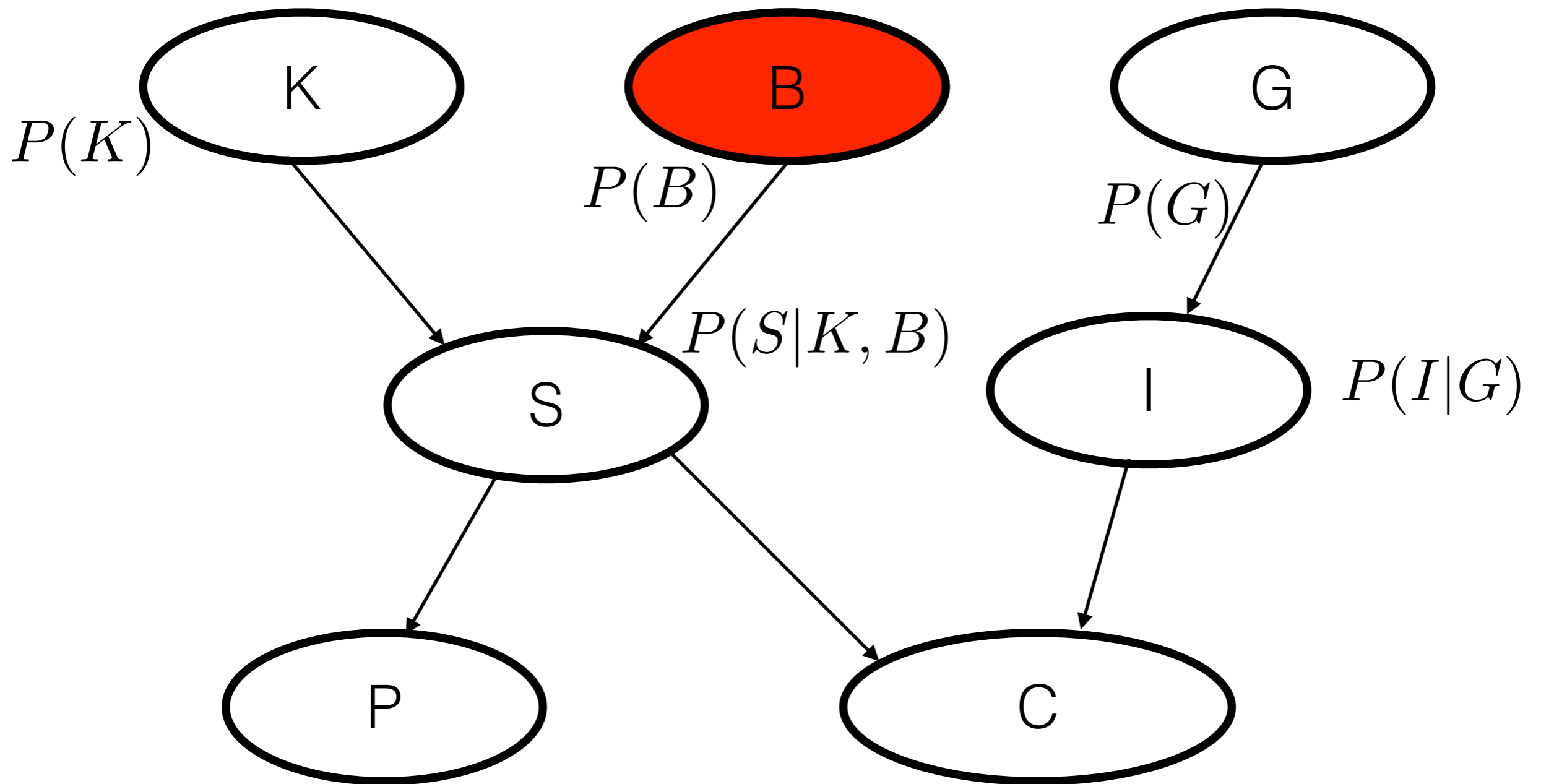
Message from node X_i to **Child** X_k on round t

$$M_{i \rightarrow k}^t(x_i) = \sum_{\text{Parents}(X_i)} E_{X_i}(x_i) P(X_i = x_i | \text{Parents}(X_i)) \text{ (product of all messages but one from } X_j \text{)} \\ \text{from previous round (t-1)}$$

After convergence:

$$P(X_i = x_i | \text{Observation}) \propto \sum_{\text{values of Parent}(X_i)} E_{X_i}(x_i) \times P(X_i = x_i | \text{Parent}(X_i)) \times \text{Product of all messages}$$

BELIEF PROPAGATION



BELIEF PROPAGATION

